

DISSERTAÇÃO DE MESTRADO
EM ENGENHARIA DE SISTEMAS ELETRÔNICOS E DE AUTOMAÇÃO

**ESTUDO SOBRE O CONSUMO DE ENERGIA EM REDES-EM-CHIP
BASEADAS EM DISPOSITIVOS NANOELETRÔNICOS**

Edylara Ribeiro Rangel

Orientadora: Janaina Gonçalves Guimarães

Brasília, agosto 2017

UNIVERSIDADE DE BRASÍLIA
FACULDADE DE TECNOLOGIA

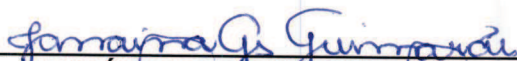
**UNIVERSIDADE DE BRASÍLIA
FACULDADE DE TECNOLOGIA
DEPARTAMENTO DE ENGENHARIA ELÉTRICA**

**ESTUDO SOBRE O CONSUMO DE ENERGIA EM REDES-EM-CHIP
BASEADAS EM DISPOSITIVOS NANOELETRÔNICOS**

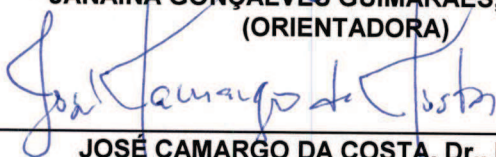
EDYLARA RIBEIRO RANGEL

DISSERTAÇÃO DE MESTRADO SUBMETIDA AO DEPARTAMENTO DE ENGENHARIA ELÉTRICA DA FACULDADE DE TECNOLOGIA DA UNIVERSIDADE DE BRASÍLIA, COMO PARTE DOS REQUISITOS NECESSÁRIOS PARA A OBTENÇÃO DO GRAU DE MESTRE.

APROVADA POR:



JANAÍNA GONÇALVES GUIMARÃES, Dra., UFSC
(ORIENTADORA)



JOSÉ CAMARGO DA COSTA, Dr., ENE/UNB
(EXAMINADOR INTERNO)



MARCELO GRANDI MANDELLI, Dr., CIC/UNB
(EXAMINADOR EXTERNO)

Brasília, 14 de agosto de 2017.

FICHA CATALOGRÁFICA

RANGEL, EDYLARA RIBEIRO

ESTUDO SOBRE O CONSUMO DE ENERGIA EM REDES-EM-CHIP BASEADAS EM DISPOSITIVOS NANOELETRÔNICOS [Distrito Federal] 2017.

xv, 70p., 210 x 297 mm (ENE/FT/UnB, Mestre, Engenharia Elétrica, 2017)

Dissertação de Mestrado - Universidade de Brasília, Faculdade de Tecnologia.

Departamento de Engenharia Elétrica

1. Redes-em-Chip

2. Interconexões

3. Nanoeletrônica

4. Energia

I. ENE/FT/UnB

II. Título (série)

REFERÊNCIA BIBLIOGRÁFICA

RANGEL, E. R. (2017). ESTUDO SOBRE O CONSUMO DE ENERGIA EM REDES-EM-CHIP BASEADAS EM DISPOSITIVOS NANOELETRÔNICOS. Dissertação de Mestrado, Publicação 673/2017 DM/PGEA.

Departamento de Engenharia Elétrica, Universidade de Brasília, Brasília, DF, 68 p.

CESSÃO DE DIREITOS

AUTOR: Edylara Ribeiro Rangel

TÍTULO: ESTUDO SOBRE O CONSUMO DE ENERGIA EM REDES-EM-CHIP BASEADAS EM DISPOSITIVOS NANOELETRÔNICOS

GRAU: Mestre em Engenharia de Sistemas Eletrônicos e de Automação ANO: 2017

É concedida à Universidade de Brasília permissão para reproduzir cópias desta dissertação de mestrado e para emprestar ou vender tais cópias somente para propósitos acadêmicos e científicos. O autor reserva outros direitos de publicação e nenhuma parte dessa dissertação de mestrado pode ser reproduzida sem autorização por escrito do autor.

Edylara Ribeiro Rangel

Depto. de Engenharia Elétrica (ENE) - FT

Universidade de Brasília (UnB)

Campus Darcy Ribeiro

CEP 70919-970 - Brasília - DF - Brasil

Dedico este trabalho a Deus, ao meu esposo, à minha família e à Profa. Janaina por acreditarem em mim.

AGRADECIMENTOS

Agradeço a Deus por me dar a oportunidade de realizar este mestrado na Universidade de Brasília, onde em um momento tão difícil, Ele me mostrou esse caminho e abriu outra porta em minha vida. Obrigada Senhor, por proporcionar esse aprimoramento do meu conhecimento e por sempre cuidar de cada detalhe em minha vida. Obrigada Senhor, pela oportunidade de entender um pouquinho mais sobre o vasto conhecimento do mundo nanométrico.

Agradeço à minha orientadora, Janaina Guimarães, por ter acreditado em mim, pelo conhecimento que compartilhou comigo, pelo respeito aos meus limites e pelo carinho. Obrigada Janaina, por dar sentido a palavra orientação, pois quando estava confusa e perdida, você conseguiu me mostrar o caminho e reavivar meu ânimo e coragem, a cada reunião. Obrigada Janaina, pelo ser excepcional que você é, que seus caminhos sejam sempre abençoados.

Agradeço ao meu esposo, Rodrigo, pois Deus o colocou na minha vida para que pudéssemos crescer juntos. Obrigada meu amor, pela sua paciência, incentivo, conselhos e por ser o meu companheiro. Essa conquista também é sua.

Agradeço aos meus pais, Aparecida e Evaluizio, que sempre incentivaram o estudo em minha vida. Obrigada mãe e pai, pelo exemplo que vocês são em minha vida.

Agradeço à minha irmã, Luisamara, pelas palavras de incentivo e apoio. Obrigada irmã, por sempre acreditar no meu potencial.

RESUMO

A evolução da indústria eletrônica que permitiu a implementação de arquiteturas de múltiplos núcleos foi motivada principalmente pelo consumo de energia, pois elas oferecem melhor desempenho e menor dissipação de potência do que os sistemas de processamento único. Com o aumento do número de núcleos em um único *chip*, a arquitetura de comunicação que interliga esses núcleos começou a ganhar importância. Assim, para resolver os problemas de interconectividade e comunicação dos sistemas em *chip*, a arquitetura de comunicação do tipo redes-em-chip (NoC - *Network-on-Chip*) tem sido proposta como uma solução altamente estruturada pela comunidade científica. Estimativas do consumo de energia das arquiteturas de comunicação devem ser realizadas no início do projeto, pois a comunicação do *chip* representa uma porção significativa do total de energia e área consumida pelo *chip*. Neste contexto, este trabalho objetiva estudar sobre o consumo de energia em NoCs baseadas em dispositivos nanoeletrônicos, por meio de um modelo analítico previamente apresentado. Para obter o consumo total de energia da comunicação do *chip*, esse modelo utiliza como base alguns parâmetros, tais como, a energia das interconexões e dos roteadores, e a distribuição de probabilidade de comunicação. O objetivo principal deste trabalho é verificar quantitativamente qual a contribuição da nanoeletrônica na redução do consumo de energia, na arquitetura de comunicação do tipo NoC, com ênfase no estudo das interconexões. Desta forma, são feitas simulações para verificar o comportamento da latência e da energia das interconexões que conectam os roteadores da rede, em função dos nós de tecnologia, bem como, é realizada a comparação do consumo de energia entre redes com roteadores nanoeletrônicos e redes com roteadores CMOS. Por fim, é realizada uma análise comparativa entre o consumo de energia de redes com interconexões de cobre e nanotubo de carbono, utilizando roteadores nanoeletrônicos. Os resultados obtidos neste trabalho mostram que a nanoeletrônica é uma tecnologia que aparenta ser uma solução promissora na redução do consumo de energia dos futuros *chips* e dispositivos.

ABSTRACT

The evolution of the electronic industry that allowed the implementation of multi-core architectures was motivated mainly by the energy consumption, since they offer better performance and less power dissipation than the single processing systems. With the increase in the number of cores on a single chip, the communication architecture that interconnects these cores began to gain importance. Thus, to solve the problems of interconnectivity and communication of the systems in chip, Networks-on-Chip (NoC) communication architecture has been proposed as a solution highly structured by the scientific community. Estimates of the energy consumption of communication architectures should be carried out at the beginning of the project because the communication of the chip represents a significant portion of the total energy and area consumed by the chip. In this context, this work aims to study energy consumption in NoCs based on nanoelectronic devices, through an analytical model previously presented. To obtain the total energy consumption of the chip communication, this model uses as base some parameters, such as the energy of the interconnections and the routers, and the Communication Probability Distribution. The main objective of this work is to verify quantitatively the contribution of nanoelectronics in the reduction of energy consumption in NoC communication architecture, with emphasis on the study of interconnections. In this way, simulations are performed to verify the latency and energy behavior of the interconnections that connect the routers of the network, as a function of the technology nodes, as well as, the comparison of the energy consumption between networks with nanoelectronic routers and networks with CMOS routers is made. Finally, a comparative analysis was performed between the energy consumption of networks with copper and carbon nanotube interconnections using nanoelectronic routers. The results obtained in this work show that nanoelectronics is a technology that appears to be a promising solution in reducing the energy consumption of future chips and devices.

SUMÁRIO

1 - INTRODUÇÃO	1
1.1 - MOTIVAÇÃO.....	3
1.2 - OBJETIVOS	4
1.3 - ORGANIZAÇÃO DA DISSERTAÇÃO.....	4
2 - FUNDAMENTAÇÃO TEÓRICA	5
2.1 - OS SISTEMAS EM CHIP E O PRINCÍPIO DAS REDES DE INTERCONEXÃO	5
2.2 - REDES-EM-CHIP - CONCEITOS BÁSICOS.....	6
2.2.1 - Topologia	8
2.2.2 - Roteadores	10
2.2.3 - Controle de Fluxo.....	11
2.2.4 - Roteamento.....	12
2.2.5 - Arbitragem	13
2.2.6 - Estratégia de Chaveamento.....	14
2.2.7 - Parâmetros de Desempenho.....	15
2.2.7.1. Latência	15
2.2.7.2. Largura de banda.....	16
2.2.7.3. Vazão	16
2.3 - INTERCONEXÕES.....	16
2.3.1 - CNT para interconexões futuras em nanoescala	18
2.3.2 - Interconexões NoC.....	19
2.3.3 - Inserção de repetidores	21
2.3.4 - Modelos de interconexão	22
2.3.4.1. Modelo de Interconexão de cobre	22
2.3.4.1.1. Resistência do cobre	23
2.3.4.1.2. Capacitância do cobre.....	24
2.3.4.1.3. Indutância do cobre.....	25
2.3.4.2. Modelo de Interconexão do SWCNT	25
2.3.4.2.1. Resistência do SWCNT isolado	26
2.3.4.2.2. Indutância do SWCNT isolado.....	27
2.3.4.2.3. Capacitância do SWCNT isolado.....	27
2.3.4.3. Modelo de Interconexão do BCNT	28
2.3.4.3.1. Resistência do BCNT	29
2.3.4.3.2. Indutância do BCNT	29
2.3.4.3.3. Capacitância do BCNT.....	29

3 - ESTIMATIVA DO CONSUMO DE ENERGIA DE REDES-EM-CHIP	31
3.1 - CONSUMO DE ENERGIA DA INTERCONEXÃO.....	31
3.1.1 - Obtenção do consumo de energia da interconexão.....	31
3.2 - CONSUMO DE ENERGIA DO ROTEADOR.....	32
3.2.1 - Arquitetura do Roteador Nanoeletrônico	32
3.2.2 - Obtenção do consumo de energia do roteador	33
3.3 - GERAÇÃO DE TRÁFEGO	34
3.3.1 - Tráfego uniforme aleatório	35
3.3.2 - Tráfego permutação de bit	35
3.3.3 - Tráfego <i>Nearest Neighbor</i>	35
3.4 - MODELANDO A LOCALIDADE ESPACIAL DE COMUNICAÇÃO NoC USANDO A REGRA DE RENT	35
3.4.1 - Geração de tráfego utilizando a regra de Rent	36
3.4.2 - Distribuição de probabilidade de comunicação da regra de Rent	37
3.5 - MODELO ANALÍTICO PARA CÁLCULO DO CONSUMO DE ENERGIA EM NoCs ...	39
4 - METODOLOGIA.....	41
4.1 - INTRODUÇÃO	41
4.2 - OBTENÇÃO DO CONSUMO DE ENERGIA DAS INTERCONEXÕES GLOBAIS DE COBRE E BCNT.....	41
4.3 - OBTENÇÃO DO CONSUMO DE ENERGIA DO ROTEADOR.....	43
4.4 - OBTENÇÃO DO CONSUMO DE ENERGIA DE UMA NOC.....	44
5 - RESULTADOS E ANÁLISES.....	46
5.1 - INTRODUÇÃO	46
5.2 - OBTENÇÃO DO CONSUMO DE ENERGIA DAS INTERCONEXÕES.....	46
5.3 - OBTENÇÃO DO CONSUMO DE ENERGIA DO ROTEADOR.....	50
5.4 - OBTENÇÃO DO CONSUMO DE ENERGIA DE NoCs.....	51
5.4.1 - Análise do Consumo de Energia de NoCs com Dispositivos Nanoeletrônicos	52
5.4.2 - Análise do Consumo de Energia de NoCs com dispositivos CMOS	53
5.4.3 - Comparação do Consumo de Energia de NoCs CMOS X NoCs Nano	54
5.4.4 - Comparação do Consumo de Energia entre NoCs com Interconexões de Cobre e BCNT.....	54
6 - CONCLUSÕES E PERSPECTIVAS FUTURAS.....	57
REFERÊNCIAS BIBLIOGRÁFICAS.....	59
APÊNDICES.....	69
A - FUNCIONAMENTO DO SET	69
B - TABELAS COMPLEMENTARES.....	69

LISTA DE FIGURAS

Figura 1-1 - Aumento do número de núcleos no tempo [3].....	2
Figura 2-1 - Estrutura básica de uma NoC.	6
Figura 2-2 - Redes diretas: (a) anel, (b) mesh, (c) toróide, d) hipercubo e (e) completamente conectada.....	9
Figura 2-3 - Redes indiretas: (a) redes multiestágio e (b) árvore-gorda.....	9
Figura 2-4 - Arquitetura de um roteador NoC.....	10
Figura 2-5 - Camadas de interconexão em processadores modernos (modificado de [33]).....	17
Figura 2-6 - Atraso relativo das interconexões em implementações ASIC [33].....	18
Figura 2-7 - Tecnologia Intel 65 nm com 8 camadas de cobre, em 2004 [45].....	19
Figura 2-8 - Interconexões NoC.....	20
Figura 2-9 - Interconexão com repetidores.....	21
Figura 2-10 - Dimensões da interconexão	22
Figura 2-11 - Modelo de interconexão de cobre.....	23
Figura 2-12 - Capacitâncias da interconexão.....	24
Figura 2-13 - Estrutura básica de um CNT. Lâmina de grafeno (esquerda), SWCNT (meio) e MWCNT (direita)[7].....	25
Figura 2-14 - Modelo de interconexão de SWCNT [50].	26
Figura 3-1 - Modelo simples de interconexão.	31
Figura 3-2 - Esquemático completo do roteador nanoeletrônico [16]	32
Figura 3-3 - Diagrama de blocos simplificado do fluxo de um flit.	33
Figura 3-4 - CPD para diferentes tipos de padrão de tráfego em uma rede 8x8 com topologia em malha. (a) Uniforme aleatório. (b) Bit complement. (c) Bit rotation. (d) Nearest Neighbor com fator de localização de 50 %. (e) Regra de Rent com expoente de 0,75.	38
Figura 4-1 - Fluxograma das etapas para obtenção do consumo de energia para as interconexões de cobre e BCNT.....	42
Figura 4-2 - Circuito utilizado na simulação das interconexões.....	43
Figura 4-3 - Fluxograma das etapas para obtenção do consumo de energia do roteador baseado em dispositivos CMOS ou nanoeletrônicos.	43
Figura 4-4 - Fluxograma das etapas para obtenção do consumo de energia de NoCs.	44
Figura 5-1 - Latência das interconexões de cobre e BCNT por nó de tecnologia para os dados disponibilizados pela INTEL.....	48
Figura 5-2 - Latência das interconexões de cobre e BCNT em função da tecnologia para os dados disponibilizados pelo ITRS.	48

Figura 5-3 - Energia por bit em função do nó de tecnologia para os materiais de cobre e BCNT, parâmetros INTEL.....	49
Figura 5-4 - Energia por bit em função do nó de tecnologia para os materiais de cobre e BCNT, parâmetros ITRS.....	49
Figura 5-5 - Comparativo de energia consumida pelo roteador Nano e pelo roteador CMOS para diferentes tamanhos de flit.	51
Figura 5-6 - Consumo total de energia de uma NoC 8X8 baseada em dispositivos Nanoeletrônicos.....	52
Figura 5-7 - Consumo total de energia de uma NoC 8X8 baseada em dispositivos CMOS.	53
Figura 5-8 - Comparativo do consumo total de energia entre uma NoC baseada em dispositivos nanoeletrônicos e uma NoC baseada em dispositivos CMOS.	54
Figura 5-9 - Comparativo do consumo total de energia entre uma NoC construída com interconexões de cobre e uma NoC construída com interconexões de BCNT, em função do tamanho da topologia da NoC.	55
Figura 5-10 - Comparativo do consumo total de energia de uma NoC construída com interconexões de cobre e outra construída com interconexões de BCNT, em função do nó de tecnologia.	56
Figura A-1 - Transistor mono-elétron.	69

LISTA DE TABELAS

Tabela 3.1 – Área e potência dos módulos do roteador nanoeletrônico [16].	34
Tabela 5.1 – Parâmetros de interconexão retirados dos dados disponibilizados pela INTEL e dos relatórios do ITRS.	47
Tabela B.1– Parâmetros obtidos a partir do modelo da interconexão global de cobre para a fonte INTEL.	69
Tabela B.2 – Parâmetros obtidos a partir do modelo da interconexão global de cobre para a fonte ITRS.	70
Tabela B.3 – Parâmetros obtidos a partir do modelo da interconexão global de BCNT para a fonte INTEL.	70
Tabela B.4 – Parâmetros obtidos a partir do modelo da interconexão global de BCNT para a fonte ITRS.	70

LISTA DE SÍMBOLOS

b	largura do canal em bits
T_{fio}	atraso de um único fio
bit_{trans}	número de bits transmitidos
n_{ciclos}	número de ciclos para que todo o tráfego seja entregue aos destinos
L_{min}	comprimento mínimo da porta do transistor
FO4	atraso de um inversor que conduz quatro inversores idênticos
L	comprimento do fio
W	largura do fio
T	espessura do fio
S	distância entre condutores em uma mesma camada
H	espessura do isolante
R_{Cu}	resistência do cobre
ρ_{Cu}	resistividade do cobre
ρ_{FS}	parâmetro de Fuchs e Sondheimer
ρ_{MS}	parâmetro de Mayadas e Shatkes
ρ_o	resistividade <i>bulk</i> do cobre
l_o	caminho médio livre dos elétrons do material de cobre
p_F	parâmetro de espalhamento de Fuchs
D	tamanho médio da região de depleção do contorno de grão,
R	coeficiente de reflexão no contorno
C_T	capacitância total do fio
C_a	capacitância de borda
C_b	capacitância de placa paralela
C_c	capacitância de acoplamento
ϵ	permissividade relativa para uma dada constante dielétrica.
μ_0	permeabilidade magnética do vácuo
L_{Cu}	indutância própria da interconexão de cobre
M_{Cu}	indutância mútua da interconexão de cobre
R_C	resistência de contato entre o metal e o nanotubo de carbono
R_q	resistência quântica do SWCNT
R_S	resistência de espalhamento do SWCNT
L_{CNT}	indutância total do SWCNT
C_q	capacitância quântica do SWCNT
C_e	capacitância eletrostática do SWCNT
h	constante de Planck
e	carga do elétron
l_{CNT}	comprimento do nanotubo de carbon
λ_{CNT}	comprimento do caminho médio livre do SWCNT
L_M	indutância magnética do SWCNT

L_K	indutância cinética do SWCNT
v_F	velocidade de Fermi
d_{CNT}	diâmetro do nanotubo de carbono
y	distância do nanotubo ao plano ligado ao terra
C_E	capacitância eletrostática do SWCNT
C_Q	capacitância quântica do SWCNT
x	distância entre os centros de nanotubos adjacentes
δ_{min}	distância de separação entre os nanotubos de carbono
n_{CNT}	quantidade de nanotubos de carbono disponível em um fio de BCNT
n_W	número de CNTs ao longo da largura do BCNT
n_T	número de CNTs ao longo da espessura do BCNT
R_{bundle}	resistência total do BCNT
L_{bundle}	indutância total do BCNT
C_{bundle}	capacitância total do BCNT
C_Q^{bundle}	capacitância quântica do BCNT
C_E^{bundle}	capacitância eletrostática do BCNT
C_{En}	capacitância entre placas paralelas próximas
C_{Ef}	capacitância entre placas paralelas afastadas de SWCNTs
C	capacitância total da carga
V	tensão de alimentação da fonte
E_{lbit}	energia consumida pela interconexão para transmitir um <i>bit</i>
E_{link}	energia consumida pela interconexão da NoC para transmitir um <i>flit</i>
$E_{roteador}$	energia total consumida pelo roteador
E_{SRWOR}	energia consumida pelo bloco registrador de deslocamento com registrador de saída
E_{DEMUX}	energia consumida pelo demultiplexador
E_{EB}	energia consumida pelo <i>buffer</i> elástico
E_{RoR}	energia consumida pelo árbitro
E_{PISO}	energia consumida pelo registrador de deslocamento
r	raio de distância utilizado pelo tráfego <i>Nearest Neighbor</i>
$P(d)$	probabilidade de comunicação entre dois núcleos
p	expoente de Rent
k	coeficiente de Rent
d	distância de Manhattan / número de saltos
Γ	coeficiente de normalização
N	número de nós na rede
$E_{flit}(d)$	energia média consumida por um <i>flit</i>
N_{flits}	número de <i>flits</i> por pacote
$N_{pacotes}$	número total de pacotes

Siglas

BCNT - *Single-Walled Carbon Nanotube Bundle*

CAD - *Computer Aided Design*

CMOS - *Complementary Metal-Oxide-Semiconductor*

CMP - *Chip Multiprocessor*

CNT - *Carbon Nanotubes*

CPD - *Communication Probability Distribution*

CPU - *Central Processing Unit*

DOR - *Dimension-Ordered*

EB - *Elastic Buffer*

FPGAs - *Field-Programmable Gate Array*

GSI - *Giga Scale Integration*

IP - *Intellectual Property*

ITRS - *International Technology Roadmap for Semiconductors*

MOSFET - *Metal Oxide Semiconductor Field Effect Transistor*

MPSoC – *Multiprocessor System-on-Chip*

MWCNT - *Multi-Walled Carbon Nanotube*

NAND - *Not AND*

nanoNAND - *NAND nanoeletrônica*

NI – *Network Interface*

NoC - *Network-on-Chip*

PISO - *Parallel-in, Serial-out Register*

RoR - *Round and Robin*

SET - *Single-Electron Transistor*

SIPO - *Serial-in, Parallel-out*

SoC - *System-on-Chip*

SRwOR - *Shift-Register with Output Register*

SWCNT - *Single-Walled Carbon Nanotube*

TSI - *Tera Scale Integration*

VLSI – *Very Large Scale Integration*

WLD - *Wire Length Distribution*

1 - INTRODUÇÃO

Com base na observação de que a densidade de componentes em circuitos integrados dobrava a intervalos regulares, o cofundador da Intel, Gordon E. Moore, declarou em seu famoso artigo de 1965 que o número de componentes em um *chip* semiconductor de menor custo, cresce exponencialmente no tempo [1]. Essa previsão ficou conhecida como Lei de Moore e continua até hoje, muito além dos 10 anos iniciais pronunciado por Moore. O crescimento da indústria eletrônica tem sido impulsionado por essa previsão que estimou que a quantidade de transistores dobraria a cada 18 meses. Essa duplicação implica em miniaturização dos componentes, corroborada pela redução das dimensões dos terminais do transistor e o aumento da velocidade de chaveamento do dispositivo.

A escala de integração de Dennard baseada no processo de miniaturização dos dispositivos, demonstrou que as dimensões e as tensões de operação do MOSFET (*Metal Oxide Semiconductor Field Effect Transistor*) devem ser reduzidas pelo mesmo fator, para manter o campo elétrico constante [2]. A Lei de Moore, em conjunto com a escala de integração de Dennard, estimulou que cada geração tecnológica produzisse o dobro do número de transistores na mesma área de silício, onde cada transistor seria 1,4 vezes mais rápido que a geração anterior, utilizando a mesma densidade de potência[3]. Essa integração de vários componentes em uma mesma pastilha de silício ficou conhecida como sistemas em *chip* (SoC - *Systems on Chip*).

A partir da integração dos transistores com dimensões inferiores a 90 nm, a corrente de fuga entre os componentes do circuito aumentou. Com isso, o limiar da tensão mínima de operação dos transistores parou de reduzir, provocando o aumento exponencial do consumo energético, a elevação da dissipação de calor e a restrição dos projetos de processadores modernos [4]. Assim, a miniaturização dos transistores foi restringida pelo problema da barreira de utilização (*utilization wall*), onde o déficit de utilização de área é superior ao previsto pela escala de Dennard, afetando diretamente o consumo de potência em circuitos integrados.

Diante do aumento exponencial da potência estática em tecnologias menores do que 90 nm, para melhorar o desempenho dos dispositivos eletrônicos e continuar com o crescimento exponencial dos dispositivos em um único *chip*, os arquitetos de computadores replicaram o número de núcleos de processamento em um *chip*, executando-os em paralelo. Com isso, surgiram os primeiros multiprocessadores em *chip* (MPSoC - *Multiprocessor System-on-Chip*) ou multicores. A Figura 1-1 destaca essa tendência traçando a evolução no tempo da quantidade de núcleos no *chip*, em várias arquiteturas comerciais e de pesquisa [3]. Como pode ser observado, o número de núcleos aumentou consideravelmente na última década.

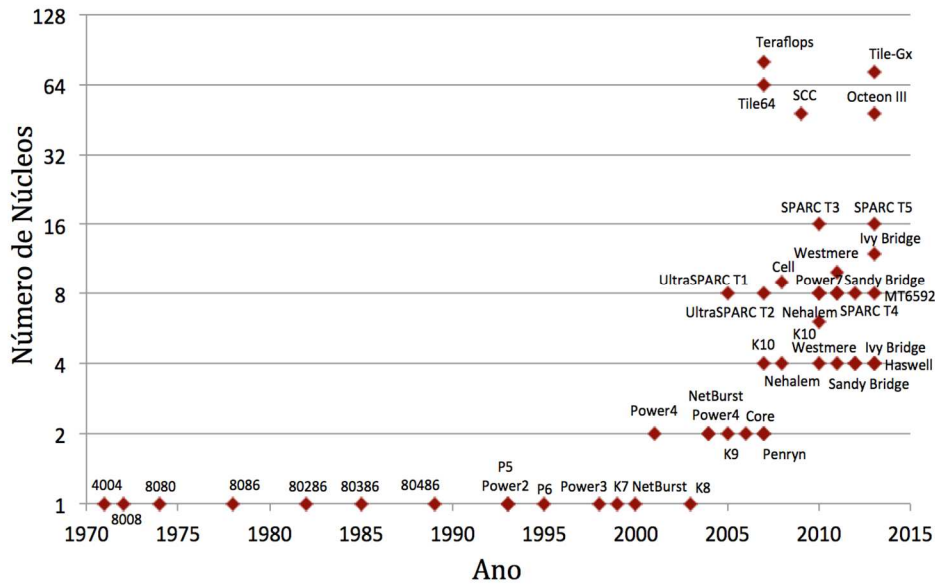


Figura 1-1 - Aumento do número de núcleos no tempo [3].

Nesse cenário, com o aumento do número de núcleos em um único *chip*, a arquitetura de comunicação que interliga esses núcleos começou a ganhar importância. Assim, para resolver os problemas de interconectividade e comunicação dos SoCs, a arquitetura de comunicação do tipo redes-em-chip (NoC - *Network-on-Chip*) tem sido proposta como uma solução altamente estruturada pela comunidade científica. Dentre os benefícios dessa arquitetura estão a grande escalabilidade e um maior nível de conectividade.

Ao contrário do que acontece com os transistores, a miniaturização das interconexões não favorece o desempenho dos circuitos. Os limites físicos das interconexões ameaçam, potencialmente, desacelerar ou até mesmo parar o progresso histórico que vem sendo alcançado pela indústria de semicondutores nos últimos anos [5]. Portanto, os limites de interconexão continuam sendo um problema na evolução dos circuitos eletrônicos integrados, especialmente em dimensões nanométricas. Novas tecnologias como as interconexões óticas, as interconexões de rádio frequência ou sem fio e as interconexões de nanotubo de carbono (CNT - *Carbon Nanotubes*) estão sendo estudadas para superar as limitações das interconexões de cobre que são as mais utilizadas atualmente [6][7][8][9].

Ainda, em escala nanométrica, os efeitos quânticos não podem ser mais desprezados e o comportamento do transistor deve se adequar à nova dimensão de operação. Dessa forma, surge a necessidade do estudo de novas tecnologias, para entendimento e desenvolvimento de novas arquiteturas. Nesse contexto, a nanoeletrônica apresenta ser uma solução promissora para continuar a redução dos dispositivos eletrônicos em escala giga (GSI - *Giga Scale Integration*) ou tera (TSI - *Tera Scale Integration*), em futuras gerações tecnológicas.

1.1 - MOTIVAÇÃO

A integração de múltiplos componentes em um único *chip* que é provocada pela miniaturização dos dispositivos, reduz o atraso de propagação da porta do transistor, resultando em frequências de operação maiores e, conseqüentemente, no aumento do consumo de potência do sistema. Além disso, com a redução do tamanho dos dispositivos, a densidade dos transistores tende a aumentar, elevando ainda mais a potência dissipada por unidade de área.

Uma determinada tarefa computacional realizada por um único processador, em uma determinada frequência, pode ser realizada por vários núcleos de processadores em paralelo, com frequência e tensão reduzidas e, deste modo, o consumo de energia é reduzido na mesma quantidade de tempo [10]. Assim, a mudança da indústria para a arquitetura de múltiplos núcleos foi motivada principalmente pelo consumo de energia, pois os MPSoCs oferecem um desempenho superior e menor dissipação de potência do que os sistemas de processamento único (CPU - *Central Processing Unit*) [11].

A questão de energia dissipada é um importante critério no projeto de SoCs. Sistemas embarcados, como *smartphones*, *tablets* e *notebooks* dependem de fonte de alimentação limitada e os processadores embutidos são projetados para minimizar o consumo de energia, a fim de aumentar a vida útil da bateria [12]. Assim, estimativas do consumo de potência das arquiteturas de comunicação devem ser realizadas no início do projeto, pois a comunicação do *chip* representa uma porção significativa do total de energia e área consumida pelo *chip* [13] [14].

As capacitâncias parasitas induzidas pelas interconexões longas, aumentam o consumo de potência dos circuitos. Esse problema é minimizado com a utilização de NoCs que utilizam interconexões ponto-a-ponto curtas entre roteadores que interconectam os elementos do *chip*. No entanto, os roteadores das NoCs também consomem potência, reduzindo a vantagem aparente em termos de consumo de potência [15]. À medida que o número de núcleos aumenta, a energia das NoCs também aumenta, impondo sérios limites de projeto no desempenho das aplicações.

Como a comunicação em *chip* consome uma parcela significativa de potência e área do *chip*, é fundamental que os roteadores sejam compactos e de baixa potência. Recentemente, um roteador desenvolvido baseado em transistores monoelétron (SET - *Single-Electron Transistor*) foi proposto para NoCs [16]. Uma breve explicação sobre o funcionamento do SET pode ser encontrada no Apêndice A. O SET ocupa uma área pequena e apresenta consumo de potência reduzido comparado aos dispositivos semicondutores de metal-óxido complementar (CMOS - *Complementary Metal-Oxide-Semiconductor*). Assim, visando reduzir a energia dissipada em uma NoC, a utilização de dispositivos nanoeletrônicos aparenta ser uma solução promissora.

1.2 -OBJETIVOS

O objetivo principal deste trabalho é verificar quantitativamente qual a contribuição da nanoeletrônica na redução do consumo de energia, na arquitetura de comunicação do tipo NoC, com ênfase na análise das interconexões. Para isso, será realizado o estudo sobre o consumo de energia em NoCs, utilizando dispositivos nanoeletrônicos baseados em SET.

Dessa forma, primeiramente será realizado o estudo do consumo de energia das partes que constituem a comunicação da NoC. Assim, o comportamento da latência e da energia das interconexões que conectam os roteadores da rede será estudado, em função da tecnologia e do material utilizado, cobre ou CNT. Em seguida, será calculado o consumo de energia dos roteadores, em função da tecnologia utilizada, CMOS ou nanoeletrônica. Após encontrar a contribuição do consumo de energia das partes que compõem a NoC, a partir do modelo analítico proposto por Bezerra [17], o consumo de energia entre redes com roteadores nanoeletrônicos e redes com roteadores CMOS será comparado. Por fim, será realizada uma análise comparativa entre o consumo de energia de redes com interconexões de cobre e CNT, ambas com roteadores nanoeletrônicos.

1.3 -ORGANIZAÇÃO DA DISSERTAÇÃO

O presente capítulo apresentou a contextualização, motivação e objetivos deste trabalho. O restante desta dissertação é resumido brevemente a seguir.

O capítulo 2 contém a fundamentação teórica dos conceitos que são necessários para a leitura desta dissertação. Assim, são apresentados os princípios de uma NoC, sua arquitetura, topologias, parâmetros e métricas utilizadas no estudo dessas redes. Ainda, será mostrada uma visão geral da estrutura de interconexões que realizam a comunicação dentro de um *chip* e serão apresentados os modelos de interconexão que serão utilizados neste trabalho.

O capítulo 3 aborda os conceitos sobre o consumo de energia das principais partes que constituem a arquitetura de comunicação do tipo NoC. Além disso, nesse capítulo é apresentado o modelo analítico usado como base neste trabalho [17], para estimar o consumo de energia da comunicação da NoC.

O capítulo 4 descreve a metodologia adotada neste trabalho, para a obtenção do consumo de energia de uma NoC, bem como, são apresentados quais os estudos realizados para atingir o objetivo principal dessa dissertação.

O capítulo 5 apresenta as conclusões obtidas com as simulações, cálculos e análises desse trabalho, e as perspectivas futuras.

2 - FUNDAMENTAÇÃO TEÓRICA

2.1 - OS SISTEMAS EM CHIP E O PRINCÍPIO DAS REDES DE INTERCONEXÃO

Os componentes de um SoC são disponibilizados em forma de módulos pré-projetados e pré-verificados, conhecidos por núcleos ou blocos de propriedade intelectual (IP - *Intellectual Property*). Um núcleo de um SoC pode ser um único processador, um módulo de memória, dispositivos de entrada e/ou saída, ou até mesmo um computador completo com processador, memória local e uma interface de rede [18]. Visando atender as exigentes demandas do mercado e a redução dos custos de projeto, é importante considerar na fabricação de um SoC, o reaproveitamento dos seus componentes, podendo estes serem desenvolvidos pela empresa responsável pelo projeto do sistema ou adquiridos de terceiros.

A comunicação em um SoC geralmente ocorre de duas maneiras: canais ponto-a-ponto ou canais multiponto. Os canais ponto-a-ponto oferecem menor latência e melhor desempenho, pois a comunicação entre dois núcleos ocorre por meio de canais dedicados, permitindo múltiplas conexões simultaneamente e proporcionando alto grau de paralelismo ao sistema [19]. Com o crescimento da quantidade de núcleos em uma única pastilha de silício, os sistemas com fios dedicados se tornam verdadeiros complexos de vias, inviabilizando a utilização dessa arquitetura de comunicação, pois esse tipo de arquitetura necessita de grande quantidade de interconexões, possuem reaproveitamento limitado e requerem projeto específico, gerando alto custo de projeto.

Para a arquitetura de canais multiponto, também conhecida por barramento, os núcleos do sistema compartilham a mesma estrutura de comunicação. Essa abordagem apresenta maior reaproveitamento, em relação à arquitetura de fios dedicados, uma vez que a mesma estrutura de comunicação pode ser reutilizada em diferentes sistemas, reduzindo tempo e custo de projeto.

Entretanto, a arquitetura de barramento possui baixo grau de paralelismo e de escalabilidade, pois como as interconexões são compartilhadas, os núcleos concorrem pelo uso do barramento e apenas um pode ser atendido, enquanto os demais esperam pela liberação do recurso. Assim, a largura de banda dessa estrutura é fixa e a adição de novos núcleos reduz o tempo de propagação dos sinais pelas interconexões do sistema, o que pode limitar seu desempenho e provocar atrasos na comunicação das aplicações. Além disso, com o aumento da quantidade de componentes no sistema, a carga capacitiva dos canais de comunicação também aumenta, ocasionando perdas e aumento da energia do sistema. Ainda, quanto à potência, o barramento exige grande quantidade de energia, tendo em vista que essa estrutura opera por difusão e cada sinal deve chegar a todos os pontos da estrutura de comunicação [19].

Nesse contexto, a arquitetura de múltiplos barramentos [20] e a hierarquia de barramentos [21], surgiram como alternativas. A arquitetura de múltiplos barramentos consiste na utilização de diversos canais multiponto, compartilhados entre os núcleos do sistema. Já a hierarquia de barramentos consiste no uso de dois ou mais canais, com características diferentes, interconectados por um circuito ponte. No entanto, a modificação proposta por essas estruturas, apenas reduz os problemas citados para a arquitetura de barramento, mas não os elimina.

O aumento da quantidade de transistores em um único *chip*, proporcionou o desenvolvimento de novas aplicações nas áreas de multimídia, telecomunicações e eletrônicos em geral. Para minimizar os problemas encontrados nas arquiteturas de comunicação apresentadas anteriormente, as redes de interconexão chaveada surgiram como uma solução quase que universal, explorando a possibilidade de aumentar o grau de paralelismo dos núcleos e o uso eficiente da largura de banda das interconexões. Assim, a solução proposta pela comunidade científica está no uso dos conceitos oriundos da área de redes de computadores aplicados no projeto de comunicação de SoCs. O uso dessas redes quando utilizadas em SoCs são denominadas de redes-em-chip (NoC - *Network-on-Chip*) [22][23].

2.2 -REDES-EM-CHIP - CONCEITOS BÁSICOS

Uma NoC é constituída basicamente de roteadores interligados por meio de interconexões. Para realizar a conexão dos núcleos de um SoC, são necessárias interfaces de rede (NI - *Network Interface*) que realizam a adaptação do protocolo dos núcleos ao da rede. Ainda, a NI também é responsável pelo serviço de comunicação dos núcleos [24]. Denomina-se nó o par roteador/núcleo. A Figura 2-1 ilustra a estrutura básica de uma NoC.

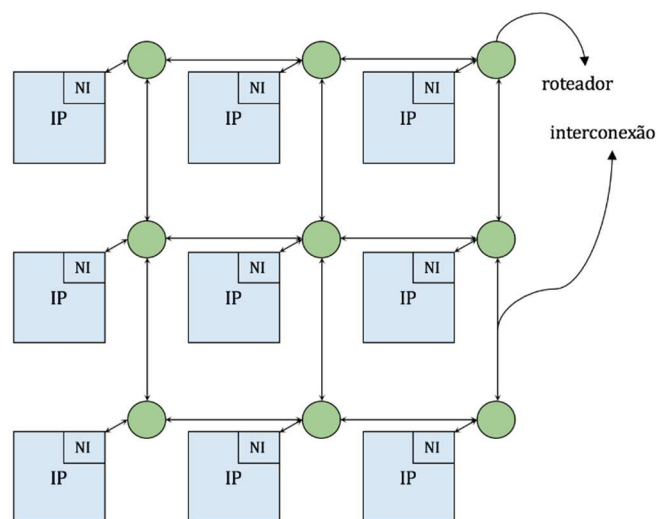


Figura 2-1 - Estrutura básica de uma NoC.

O roteador é um componente básico que transfere informações entre os núcleos, conectando um número de canais de entrada, a um número de canais de saída. Já os enlaces ou interconexões realizam a interligação de um roteador a outro roteador ou a um outro nó do sistema. Neste trabalho utiliza-se o termo salto para referenciar a conexão física entre dois núcleos/roteadores vizinhos. As interconexões podem ser unidirecionais ou bidirecionais. Geralmente, as interconexões mais utilizadas em redes de interconexão são as bidirecionais, pois permitem a transferência simultânea de informação nas duas direções do enlace.

Os dados são transferidos sob forma de mensagens, as quais são divididas em unidades menores conhecidas por pacotes que contêm palavras com tamanho igual à largura física do canal, sendo denominadas de *phit* (*physical unit*). O pacote geralmente possui estrutura semelhante à da mensagem, constituindo-se de um cabeçalho (*header*), dados úteis (*payload*) e *trailer*. As informações contidas no cabeçalho do pacote definem o caminho a ser percorrido pela mensagem, enquanto que o *trailer* é utilizado para verificação de erros e sinalização de fim do pacote. Uma mensagem é transmitida entre os núcleos perfazendo vários saltos entre os núcleos de origem e destino.

À medida que a complexidade e a integração de um SoC continuam a aumentar, muitos projetistas de sistema estão preferindo rotear pacotes, não fios [23]. As vantagens de uma NoC foram resumidas como segue [23][17]:

- Parâmetros elétricos previsíveis: Os fios não estruturados têm capacitâncias parasitas e ruídos de diafonia que são difíceis de prever. Como resultado, para garantir a confiabilidade, circuitos conservadores devem ser usados para conduzir e receber esses fios, levando a um consumo de energia excessivo. Os fios bem estruturados e previsíveis de uma NoC permitem circuitos de alto desempenho, que podem reduzir a dissipação de energia e aumentar a propagação do fio, ao mesmo tempo em que melhoram a largura de banda;
- Interface universal: ao introduzir uma interface universal para os núcleos, os componentes podem ser reutilizados em muitos sistemas, reduzindo assim a complexidade e simplificando a implementação de circuitos;
- Reusabilidade: possibilidade de aproveitar a mesma arquitetura de comunicação em aplicações distintas;
- Fator de serviço (*Duty factor*) das interconexões é otimizado: em arquiteturas de interconexões dedicadas, apenas 10% dos fios permanecem ativos em relação ao tempo total de processamento da aplicação. O fluxo agregado de informações em NoCs de propósito geral pode fornecer fator de serviço dos fios próximo a 100%;

- Permitir o uso de estratégias tolerantes a falhas: com o dimensionamento da tecnologia e a diminuição do uso da tensão, os fios tornam-se mais suscetíveis ao ruído e às falhas. Eventualmente, será impossível evitar completamente esses erros (chamados de distúrbios) na comunicação e o sistema deve ser capaz de lidar com eles. Uma arquitetura NoC pode implementar protocolos de identificação e correção de erros que tornam o sistema tolerante a falhas;
- Paralelismo (*pipelining*) das interconexões: globalmente, os protocolos assíncronos permitem o paralelismo das interconexões, aumentando assim a largura de banda, devido à multiplicidade de caminhos possíveis em uma NoC;
- Escalabilidade: a arquitetura NoC é escalável. Isso significa que com o acréscimo de um componente na rede, o número de canais de comunicação aumenta e conseqüentemente a largura de banda agregada aumenta com o tamanho da rede.

Uma rede de interconexão é caracterizada por sua topologia e por um conjunto de protocolos que definem a forma como ocorrerá a transferência de dados pela rede. No projeto de uma NoC, faz-se necessário escolher apropriadamente os diferentes requisitos de rede, pois a escolha desses parâmetros interfere diretamente no desempenho da aplicação. Assim, entender os princípios básicos das redes é de suma importância e seu desempenho pode ser avaliado por meio de algumas métricas, como largura de banda, vazão e latência [19]. As características principais de uma NoC são resumidas a seguir e uma descrição mais aprofundada pode ser encontrada em [24][25].

2.2.1 - Topologia

A topologia define a organização física da rede composta pelos nós. Ou seja, a topologia define os caminhos possíveis entre todos os nós. As topologias utilizadas em NoCs podem ser agrupadas em dois grandes grupos: redes diretas e redes indiretas.

Nas redes diretas, cada roteador está associado diretamente a um núcleo. Em termos de conectividade, a rede direta ideal é aquela que está completamente conectada, onde cada nó está interligado a todos os outros da rede (Figura 2-2 e). Porém, sua escalabilidade é limitada, pois para uma grande quantidade de núcleos, seu custo é altamente elevado. Assim, como alternativas para esse tipo de topologia, outras soluções foram propostas, tais como, as redes em anel, *mesh*, toróide (*torus*) e hipercubo. A Figura 2-2 mostra as redes diretas citadas que são as principais encontradas na literatura.

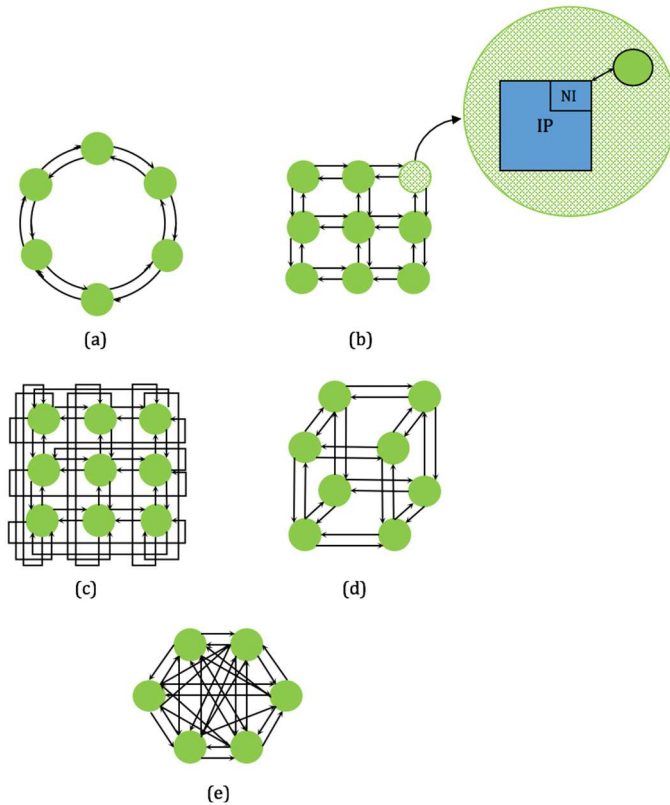


Figura 2-2 - Redes diretas: (a) anel, (b) *mesh*, (c) toróide, (d) hipercubo e (e) completamente conectada

Nas redes indiretas, somente alguns roteadores possuem ligação com núcleos. Neste tipo de rede, apenas os núcleos terminais são conectados aos roteadores e cada roteador pode conectar outros roteadores e/ou nós terminais. Entre as redes indiretas encontradas na literatura, destacam-se as redes multiestágio e árvore-gorda, apresentadas na Figura 2-3. As FPGAs (*Field-Programmable Gate Array*) são exemplos de sistemas que utilizam redes indiretas.

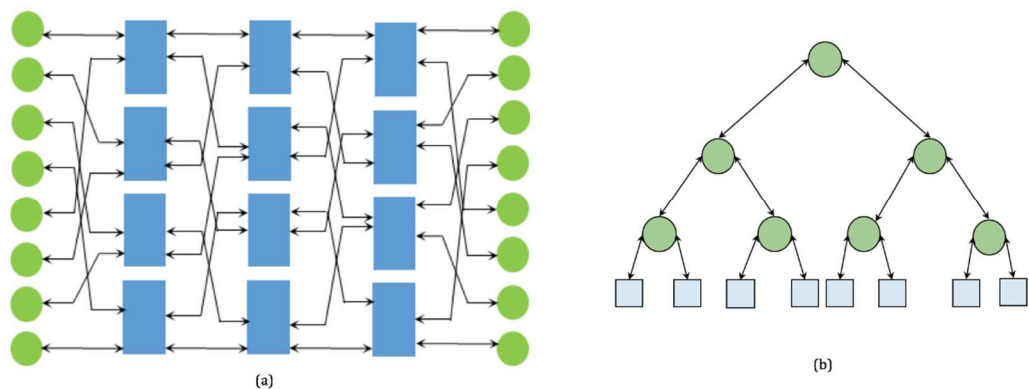


Figura 2-3 - Redes indiretas: (a) redes multiestágio e (b) árvore-gorda.

Uma das principais propriedades da topologia é a ampliação da bissecção da largura de banda. A bissecção da largura de banda é o número de fios que devem ser

cortados quando a rede é dividida em dois conjuntos iguais de nós. À medida que mais núcleos são conectados à rede, maior é o volume de comunicação e mais largura de banda é necessária. Se a largura de banda não escala adequadamente com o número de núcleos, o tráfego excessivo elevará a latência da mensagem e o desempenho do sistema será reduzido. No entanto, redes com grande bisseção de largura de banda vão exigir mais roteadores e mais fios por núcleo que consomem área considerável e aumentam o custo do sistema [17].

2.2.2 - Roteadores

Os roteadores são responsáveis por estabelecer o caminho por onde serão transferidos os dados da rede. Um roteador é normalmente constituído por um núcleo de chaveamento (*crossbar*), uma lógica de controle para roteamento e arbitragem, portas de entrada e saída para comunicação com outros roteadores e *buffers*. Ainda, as portas podem possuir controladores de enlace para implementação do protocolo físico de comunicação [19]. A Figura 2-4 mostra a arquitetura básica de um roteador NoC de cinco portas, para uma rede *mesh*.

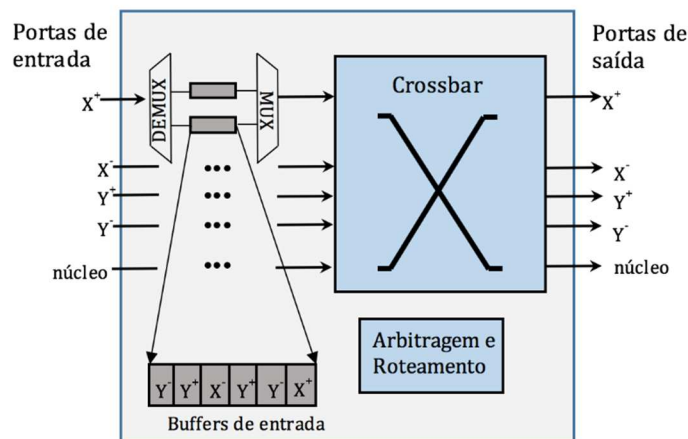


Figura 2-4 - Arquitetura de um roteador NoC.

Os *buffers* são utilizados para armazenar temporariamente os dados que não podem ser roteados imediatamente. Entretanto, em NoCs, eles possuem alto custo quanto ao consumo de energia [22]. Ainda, os *buffers* podem ser associados às portas de entrada e saída do roteador. Porém, para economizar área e energia, normalmente eles não são implementados nas portas de saída [26].

O *crossbar* é o elemento de chaveamento responsável por conectar todas as entradas do roteador à todas suas saídas, direcionando os dados de entrada à porta de saída definida pelo algoritmo de roteamento.

A unidade lógica de controle para roteamento e arbitragem é a responsável por decodificar o destino da mensagem de entrada e calcular a porta de saída mais adequada para a transmissão da mensagem, com base no algoritmo de roteamento. Ainda, essa unidade leva em conta os problemas de chaveamento e controle de fluxo que serão explicados adiante.

2.2.3 - Controle de Fluxo

A transmissão de um *flit*, ou unidade de controle de fluxo (*flow control unit*), entre as portas de entrada e saída em um roteador, é uma tarefa executada pela técnica de chaveamento. No entanto, o controle de fluxo é responsável pela administração do avanço da informação entre os roteadores. Os *buffers* são recursos temporários para armazenamento de *flits*, mas não são infinitos. As técnicas de controle de fluxo são responsáveis por determinar quando os *flits* podem ser encaminhados avaliando a capacidade dos *buffers* e a largura de banda do canal [26], garantindo que os recursos de rede (*buffers* e canais) não estejam inativos quando houver *flits* esperando para usá-los. Enquanto a topologia e o algoritmo de roteamento definem as características teóricas de latência e taxa de transferência para um determinado padrão de tráfego, o controle de fluxo é quem determina o quão perto dessa capacidade teórica a rede pode operar [3]. Existem três mecanismos principais de controle de fluxo que são comumente usados: *ack/nack*, *Stop and Go* e baseado em créditos.

O mecanismo de controle de fluxo *ack/nack* ou *handshake* é baseado em reconhecimento de dados e não faz controle do espaço do *buffer* do nó receptor. Quando um *flit* chega a um *buffer*, se o *buffer* tiver espaço disponível, o *flit* é aceito e um sinal de confirmação (*ack*) é enviado de volta ao roteador de origem. Caso não haja espaço disponível, o *flit* é descartado e um reconhecimento negativo é enviado (*nack*). Neste tipo de controle, o *flit* deve ser retido na sua origem até receber um reconhecimento positivo [26]. Assim, o transmissor realiza o reenvio do *flit* até receber um sinal *ack* como resposta. A implementação desse controle é simples, porém, a retransmissão dos dados descartados ocasiona uso ineficiente da largura de banda da rede.

O controle de fluxo *Stop and Go* surgiu como uma alternativa para reduzir a sinalização (tráfego de controle) entre o nó emissor e o nó receptor. No controle de fluxo *Stop and Go*, um *bit* de controle é mantido no *buffer* do nó receptor que indica se o nó pode ou não receber dados. O controle *Stop and Go* possui dois limites de controle correspondentes a determinados tamanhos de latência calculados, a partir do tempo de ida e volta do percurso. Quando o espaço ocupado no *buffer* atinge o limite máximo, um sinal *stop* é enviado de volta para o nó emissor, indicando que o *buffer* do receptor está cheio e solicitando parar a transmissão. Esse limite leva em conta se há espaço suficiente no *buffer* do receptor, para os *flits* que ainda estão sendo transmitidos pelo

emissor. Quando a ocupação do *buffer* de destino diminui abaixo ou igual ao limite mínimo, um sinal *go* é enviado para reativar o fluxo de *flits* [26].

No controle de fluxo baseado em créditos, cada emissor mantém uma contagem de créditos que é igual ao número de *flits* que ainda podem ser armazenados no *buffer* do lado do receptor. Sempre que um *flit* é encaminhado para o *buffer* do receptor, como ele ocupa um espaço no *buffer* de destino, o contador é decrementado. Quando o contador é zerado, isso significa que não há espaço de *buffer* disponível na outra extremidade e nenhum *flit* pode ser encaminhado. Por outro lado, sempre que um *flit* é encaminhado e libera espaço de armazenamento no *buffer* de destino, um crédito é enviado de volta ao roteador de origem, para incrementar o contador. A desvantagem deste mecanismo de controle de fluxo é a quantidade significativa de sinalização de crédito enviada para trás, o que poderia afetar o desempenho da rede [26].

Em uma rede sem congestionamento, se o controle de fluxo adotado for o *ack/nack*, espera-se que sejam necessários dois ciclos para se transmitir um *flit*, enquanto que para o controle de fluxo baseado em créditos, apenas um ciclo de relógio é necessário para transmissão de um *flit* [31].

2.2.4 - Roteamento

Os algoritmos de roteamento são executados pelos roteadores, com o objetivo de decidir qual trajeto será utilizado dentro de uma topologia, para levar uma mensagem do seu nó de origem até o seu nó de destino. A escolha do algoritmo é de extrema importância, pois afeta diretamente o desempenho da rede.

Os algoritmos de roteamento devem evitar três principais problemas: *deadlock*, *livelock* e *starvation*. O *deadlock* ocorre quando existe uma dependência cíclica entre os roteadores, onde há um impasse pela requisição de determinado serviço e nenhum deles consegue avançar. O *livelock* ocorre quando um pacote se mantém trafegando permanentemente na rede porque os canais necessários para ele chegar ao nó de destino, nunca se encontram livres. O *starvation* ocorre quando uma porta requisita um recurso da rede, mas nunca é atendida, pois possui baixa prioridade.

Os algoritmos de roteamento podem ser classificados em determinísticos ou adaptativos. No caso determinístico, o caminho escolhido para um determinado par de origem e destino é sempre o mesmo. Já no caso do adaptativo, o algoritmo utiliza alguma informação de tráfego da rede, para escolher o caminho que será percorrido por um determinado pacote. Este trabalho não aprofunda sobre os algoritmos de roteamento e apenas introduz o mais comum encontrado na literatura para NoCs, conhecido como algoritmo de roteamento XY, o qual foi adotado nessa dissertação.

O algoritmo de roteamento ortogonal XY, também conhecido como DOR (*Dimension-ORdered*), é um algoritmo determinístico, livre de *deadlock*. Neste algoritmo, a mensagem é encaminhada em uma ordem de dimensão estabelecida. Por exemplo, em uma topologia *mesh* 2D, a mensagem é encaminhada até atingir a abscissa X do nó de destino. Após chegar na ordenada X, mudanças de direção quando um pacote está na direção Y são proibidas e a mensagem segue em direção a coordenada Y do nó de destino, onde terminará sua trajetória.

2.2.5 - Arbitragem

Como visto anteriormente, um roteador é composto por várias portas de entrada e saída com seus *buffers* e canais associados. As requisições das portas de entrada do roteador, de acordo com decisões de roteamento, podem solicitar a mesma porta de saída [26]. Nesse cenário, uma operação de arbitragem é necessária para resolução de conflitos decorrentes de múltiplos pacotes competindo pela mesma porta de saída. Assim, enquanto o roteamento é um mecanismo de seleção de saída, a arbitragem é um mecanismo de seleção de entrada [19].

A operação de arbitragem introduz um atraso para determinar a atribuição das diferentes portas de saída. Essas operações são críticas em um ambiente de NoCs, devendo ser executadas com rapidez, para manter baixa a latência do sistema. O objetivo principal de um mecanismo de arbitragem é proporcionar equidade entre todas as portas de entrada, ao mesmo tempo em que obtém correspondências máximas entre solicitações e recursos. Embora existam muitas propostas para algoritmos de arbitragem e implementações, as duas principais técnicas de arbitragem para atribuir prioridades entre as requisições de entrada são: prioridade fixa e *round-robin* [26].

Na arbitragem por prioridade fixa, o árbitro atribui uma prioridade estática às requisições de cada porta de entrada. Neste mecanismo, a arbitragem é simples, mas se um dos *buffers* de entrada com maior prioridade continuar solicitando a saída associada, as entradas com menor prioridade ficam bloqueadas, onde uma requisição com menor prioridade pode nunca ser atendida e ocasionar o problema de *starvation* [26].

Na arbitragem *round-robin*, o árbitro implementa ciclos de prioridades entre todas as portas de entrada, atribuindo a prioridade mais baixa para a requisição da porta de entrada cujo pedido foi atendido pela última vez. Esta técnica de arbitragem introduz equidade entre os solicitantes, mas é mais complexa de ser implementada [26].

2.2.6 - Estratégia de Chaveamento

O chaveamento define como a mensagem será transferida dentro da rede. Os dois métodos de transferência de pacotes utilizados são: chaveamento por circuito e chaveamento por pacote.

No chaveamento por circuito (*circuit switching*), o percurso de uma mensagem do seu núcleo fonte até o núcleo de destino é estabelecido antes da comunicação ser realizada [27]. Dessa forma, os canais físicos são reservados durante a transmissão da mensagem, não podendo ser usados por outros nós da rede. Para estabelecer a rota, uma quantidade mínima de *buffers* é utilizada para armazenar os cabeçalhos que irão reservar a rota. As vantagens desse método são a latência garantida [27] e a facilidade de implementação [24], podendo ser bastante vantajoso quando as mensagens são frequentes e longas. A desvantagem é que os canais não são compartilhados. Assim, enquanto os canais estiverem reservados para um determinado fluxo, nenhuma outra mensagem poderá ser transmitida, mesmo que a conexão estiver ociosa e, com isso, os canais são subutilizados, gerando menores taxas de transferência da rede.

O chaveamento por pacote (*packet switching*) é uma técnica mais granular, onde uma mensagem é dividida em vários pacotes e os canais são reservados apenas durante a transmissão de um único pacote. Cada pacote contém um cabeçalho com as informações necessárias para sua transmissão e um número de sequência para remontagem da mensagem, após todos os pacotes chegarem ao destino. Essa técnica necessita de mais *buffers* para armazenamento dos pacotes, até a montagem da mensagem, porém, permite melhor utilização da rede, pois não reserva recursos. O chaveamento por pacote pode ser classificado em: chaveamento por pacote *Store-and-Forward*, chaveamento por pacote *Virtual-Cut-Through* e chaveamento por pacote *Wormhole*.

Na técnica *store-and-forward* (armazenar e passar), quando um pacote chega a um roteador, ele armazena completamente o pacote em seu *buffer*, antes de decidir qual será a porta de saída e o percurso para a transmissão do pacote. Portanto, os *buffers* das portas de entrada do roteador devem ser grandes o suficiente para armazenar um pacote. Assim, essa técnica possui maiores requisitos de *buffer* do que a técnica de chaveamento por circuitos. Além disso, a latência da comunicação será proporcional ao tamanho do pacote, multiplicada pelo número de saltos do percurso total.

Como alternativa ao método *Store-and-Forward*, o chaveamento por pacote *Virtual-Cut-Through* foi proposto. A diferença básica está no armazenamento do pacote, sendo este roteado no instante em que o recurso para o próximo nó estiver disponível. Com isso, a latência da comunicação da rede é reduzida. Ainda, os *buffers* deverão ser dimensionados para conter um pacote inteiro, para os casos onde a rede

estiver completamente carregada. Assim, o chaveamento *Virtual-Cut-Through* se comportará como o chaveamento *Store-and-Forward* quando o canal estiver ocupado.

O chaveamento por pacote *Wormhole* [28] é uma variação da técnica *Virtual-Cut-Through*, onde os *buffers* dos roteadores são dimensionados para armazenar apenas alguns *flits*. O *flit* é a menor unidade de informação que pode ser transmitida por meio de um canal [29]. No chaveamento *Wormhole* apenas o *flit* de cabeçalho contém as informações de roteamento. Assim, o *flit* cabeçalho é o responsável por estabelecer o percurso de todos os *flits* restantes do pacote, sendo o canal reservado até o término da transmissão do pacote. Os roteadores *Wormhole* aumentam a eficiência da comunicação entre os núcleos da rede, por meio da redução da latência dos pacotes e aumento da taxa de transferência. A principal vantagem dos roteadores *Wormhole* é a baixa necessidade de requisitos de *buffer*. No entanto, na trajetória do *flit* cabeçalho, o recurso deve estar disponível, com espaço para armazenamento em *buffer* no *crossbar* do roteador e um *flit* de largura de banda. Caso contrário, o pacote não conseguirá avançar pela rede, bloqueando o canal em utilização [19].

Para superar o problema de contenção induzido pelo chaveamento *Wormhole*, a utilização de canais virtuais foi proposta [30]. Os canais virtuais permitem a divisão dos *buffers* de entrada do roteador, permitindo que o canal físico seja compartilhado por diversos pacotes, aumentando a taxa de transferência da rede e reduzindo a sua latência [26].

2.2.7 - Parâmetros de Desempenho

A avaliação de desempenho de uma NoC tem por objetivo verificar e avaliar os serviços de uma rede. Nesta seção serão mostrados os parâmetros de desempenho de latência, vazão e largura de banda utilizados no projeto de uma NoC.

2.2.7.1. Latência

A latência de rede é o tempo decorrido a partir do momento em que o cabeçalho do pacote é injetado na rede, até o momento em que o último *flit* do pacote é recebido no nó de destino [25]. A latência é medida em unidades de tempo. Porém, como muitas comparações são realizadas utilizando-se simuladores de rede, a latência pode ser medida em ciclos de relógio gastos para o pacote percorrer um caminho [19].

Na ausência de contenção, a latência do *flit* é definida pela soma de dois fatores determinados pela topologia, o atraso do roteador e o atraso da interconexão. O atraso do roteador é tempo gasto pelo roteador para processar um único *flit*, enquanto que o atraso da interconexão é o tempo gasto para transmitir o *flit*, no fio que interliga dois roteadores. Assim, a latência é um problema crítico de projeto em vários sistemas, como por exemplo em sistemas em tempo real.

2.2.7.2. Largura de banda

A largura de banda (*bandwidth*) de um canal é a taxa máxima para transmissão de dados, em uma rede de interconexão. A Equação (2.1) apresenta o cálculo para a largura de banda [31], onde b é a largura do canal em bits, ou seja, o número total de fios da interconexão e T_{fio} é o atraso de um único fio. Assim, a unidade de medida utilizada para mensurar a largura de banda de uma rede é o “bit por segundo” (ou bps).

$$largura\ de\ banda = \frac{b}{T_{fio}} \quad (2.1)$$

2.2.7.3. Vazão

A vazão (*throughput*) em determinado canal é a quantidade máxima de dados transmitidos durante determinado intervalo de tempo, sendo expressa pela Equação (2.2), onde bit_{trans} é o número de bits transmitidos e n_{ciclos} é o número de ciclos para que todo o tráfego seja entregue aos destinos. O parâmetro vazão é medido através da contagem do número de bits que chegam no destino em um intervalo de tempo para cada fluxo (par origem-destino) [31].

$$vazao = \frac{bit_{trans}}{n_{ciclos} \cdot T_{fio}} \quad (2.2)$$

2.3 -INTERCONEXÕES

A redução dos dispositivos eletrônicos em escala nanométrica, aumentou consideravelmente o número de transistores em um *chip*, assim como o número de interconexões. À medida que as dimensões dos dispositivos diminuem e o poder de processamento dos dispositivos continua a melhorar, os atrasos de interconexão global começam a dominar os atrasos de porta dos transistores. Com isso, a comunicação entre processadores se torna um gargalo no desempenho do circuito [32].

Os primeiros processadores CMOS foram fabricados com uma única camada de metal. Para acomodar o grande número de fios em um mesmo *chip*, camadas de metal extras foram adicionadas. Por muitos anos, apenas duas ou três camadas de metal foram utilizadas para interligar os vários componentes do *chip*. Com o avanço da tecnologia no polimento químico-mecânico e em alguns outros processos semicondutores, o número de camadas de metal tem aumentado acima da camada ativa de silício [32]. Como pode ser observado na Figura 2-5, as camadas inferiores das

interconexões são mais finas e estreitas, enquanto que as camadas superiores são mais grossas e largas.

As camadas metálicas são divididas em quatro grupos distintos: locais, intermediárias, semiglobais e globais. A Figura 2-5 mostra a ideia dessa divisão que é feita com base no comprimento das interconexões. As interconexões locais geralmente conectam os componentes de um dispositivo eletrônico, enquanto que as interconexões intermediárias e semiglobais são utilizadas para conectar os dispositivos dentro de um mesmo bloco funcional. Já as interconexões globais são utilizadas para interligar os componentes distantes no *chip*, além de serem responsáveis pela distribuição de energia e sinal de *clock*. A maior seção transversal destes fios, decorrente de sua maior largura e espessura, garante menor resistência e, portanto, aumento da velocidade de propagação. Portanto, a arquitetura de rede de interconexão multinível não é apenas um requisito para o roteamento, mas também uma solução parcial para o problema de latência da interconexão [35].

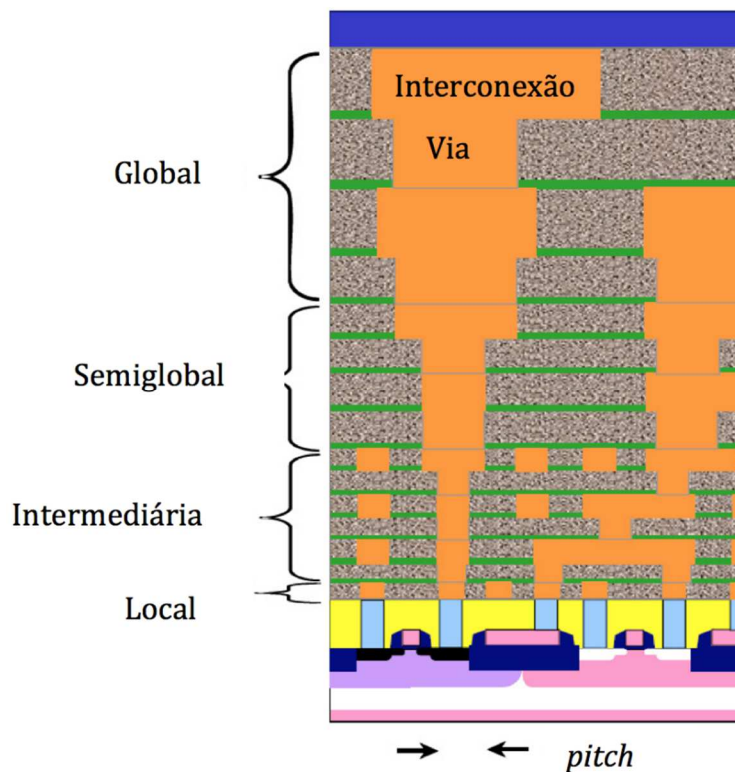


Figura 2-5 - Camadas de interconexão em processadores modernos (modificado de [33]).

Porém, com a redução da dimensão dos dispositivos eletrônicos, a largura e a espessura das interconexões são reduzidas, inclusive as globais. Como resultado, a resistência aumenta e, à medida que as interconexões se aproximam, a capacitância de acoplamento entre os fios adjacentes também aumenta, elevando o atraso resistência-capacitância (RC) [36]. A Figura 2-6 mostra o atraso relativo dos fios *versus* a redução

do tamanho da tecnologia. As duas curvas superiores do gráfico destacam o aumento alarmante no atraso das interconexões globais, em relação aos atrasos de porta dos transistores, à medida que os tamanhos das dimensões dos dispositivos diminuem para escala nanométrica. Os atrasos das interconexões globais aumentam exponencialmente com a redução do tamanho da tecnologia de processamento ou, na melhor das hipóteses, linearmente após a inserção de repetidores [37].

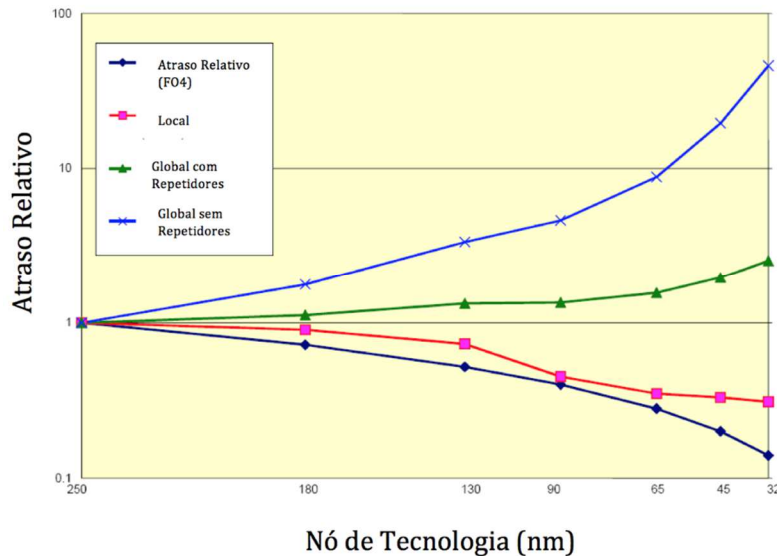


Figura 2-6 - Atraso relativo das interconexões em implementações ASIC [33].

2.3.1 - CNT para interconexões futuras em nanoescala

No passado, o alumínio (Al) foi substituído pelo Cobre (Cu) para melhorar o produto do atraso RC das interconexões, pois o cobre oferece maior condutividade em comparação ao alumínio [39] e possui uma maior resistência à eletromigração, em alta densidade de corrente [40]. Em comparação com o alumínio, o cobre pode suportar cerca de cinco vezes mais densidade de corrente para aplicações em circuitos integrados [41].

Com o avanço da integração da tecnologia em escala muito grande (VLSI – *Very Large Scale Integration*), o número de interconexões no *chip* aumentou. Para acomodar essa maior quantidade de interconexões, as dimensões da seção transversal do fio reduziram na ordem do caminho médio livre dos elétrons do cobre (aproximadamente 40 nm em temperatura ambiente) [42]. A redução das dimensões das interconexões conduz a um aumento significativo da resistividade do cobre, a cada geração tecnológica. Esse aumento é ocasionado devido à eletromigração e ao aumento do espalhamento de contorno e superfície [44].

A capacitância associada às interconexões determina diretamente tanto o atraso RC de interconexão, quanto a dissipação de energia dinâmica de interconexão.

No esforço de reduzir essa capacitância, progressivamente materiais dielétricos com baixa permissividade elétrica (baixo-k) foram introduzidos em muitas gerações tecnológicas [40]. A Figura 2-7 mostra uma imagem de seção transversal, realizada com microscópio de varredura (SEM - *Scanning Electron Microscope*), que apresenta a estrutura de interconexão do cobre na tecnologia de 65 nm.

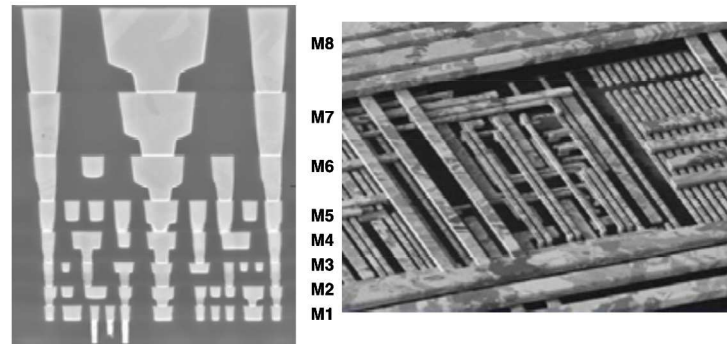


Figura 2-7 - Tecnologia Intel 65 nm com 8 camadas de cobre, em 2004 [45].

As interconexões de cobre estão passando por problemas similares aos encontrados nas interconexões de alumínio. Assim, os CNT metálicos são vistos como um potencial substituto para as interconexões de cobre, devido as suas propriedades elétricas, térmicas e mecânicas [46][47]. A alta condutividade térmica do CNT, permite suportar densidades de até 10^{14} A/m², sendo que o cobre suporta densidades inferiores a 10^{11} A/m². Essas propriedades possibilitam uma tolerância superior do CNT à eletromigração, em comparação ao cobre. Como a eletromigração causa problemas de confiabilidade a longo prazo, os CNT metálicos aparentam ser a melhor opção para interconexões futuras em larga escala [6][7].

2.3.2 - Interconexões NoC

Em um projeto de NoC, as interconexões desempenham um papel fundamental na comunicação e podem gerar um grande impacto no consumo total de energia, na área de fiação e no desempenho do sistema. Um dos desafios mais críticos de um projeto NoC é fornecer a largura de banda necessária estabelecida pelo projeto SoC que visa alcançar determinado limiar de desempenho. Como a tecnologia está encaminhando para o domínio nanométrico, alcançar maior largura de banda para os canais de comunicação, se torna uma tarefa cada vez mais desafiadora.

As interconexões da NoC são constituídas por fios de sinal em paralelo, com largura e espaçamento fixos, conforme mostrado na Figura 2-8. Conforme visto anteriormente, o desempenho do sistema pode ser mensurado utilizando a métrica de largura de banda que é a mais utilizada na literatura.

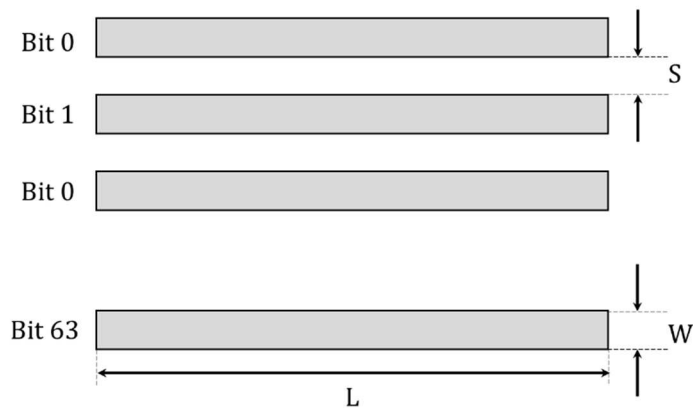


Figura 2-8 - Interconexões NoC.

Para obter maior largura de banda, é importante que o atraso seja mínimo. No projeto de um NoC, o tempo de ciclo de *clock* mínimo concebível pode ser assumido para ser igual ao valor de 15 FO4, onde FO4 é definido como o atraso de um inversor que conduz quatro inversores idênticos [38]. Em diferentes nós de tecnologia, o FO4 pode ser estimado pela Equação (2.3), onde L_{min} é o comprimento mínimo da porta em qualquer nó de tecnologia [34]. Na Equação (2.3), L_{min} deve ser inserido na escala em micro.

$$FO4 = 425 \cdot L_{min} \quad (2.3)$$

Em longos fios, o atraso intrínseco de RC pode exceder facilmente o limite de 15 FO4, o que pode limitar o tempo de ciclo de *clock* do projeto e, como consequência, reduzir a largura de banda do sistema. O comprimento das interconexões de uma NoC é definido em função da topologia, do número de núcleos e do tamanho do *chip*. Dessa forma, dependendo do comprimento dos fios, diferentes técnicas podem ser necessárias para reduzir o atraso RC intrínseco [34].

Na NoC, as interconexões que interligam os roteadores são as mais longas, atrás somente daquelas responsáveis por entregar a energia, *clock* e terra ao sistema [48]. As interconexões globais são adequadas e recomendadas para as interconexões NoC, sendo apropriadas para conseguir o fornecimento de largura de banda de um projeto SoC [34]. É extremamente importante que os SoCs sejam projetados com grande quantidade de largura de banda, para satisfazer a alta densidade de comunicação entre os processadores, pois quanto maior a largura de banda, maior a quantidade de dados transmitidos e menor a contenção de informações dentro da rede.

Para obter maior largura de banda em uma NoC, é possível projetar roteadores

pipelined de tal forma que processem um *flit* por ciclo. Porém, a duração do ciclo de *clock*, geralmente determina a velocidade com que cada *flit* pode ser processado na rede. Em redes nanométricas, o ciclo de *clock* não é limitado pela frequência de operação dos dispositivos, mas pelas ligações entre dois roteadores [34]. Neste contexto, a arquitetura de comunicação que interliga os núcleos dentro de um processador será um dos fatores limitantes que determinará o desempenho do sistema.

2.3.3 - Inserção de repetidores

A resistência e capacitância de um fio são funções lineares do comprimento do fio. Assim, o atraso de propagação do fio interrupto é uma função quadrática do comprimento do fio. Para interconexões NoC mais longas, o dimensionamento e o espaçamento do fio isoladamente podem não ser suficientes para limitar esse crescimento quadrático. A adição de repetidores, tais como, inversores e *buffers*, a intervalos regulares ao longo do fio é uma técnica padrão para reduzir o atraso da interconexão. Nesta técnica, a divisão do fio em múltiplos segmentos faz com que o atraso do fio seja uma função linear do comprimento do fio. A Figura 2-9 mostra um fio de interconexão com repetidores inseridos.

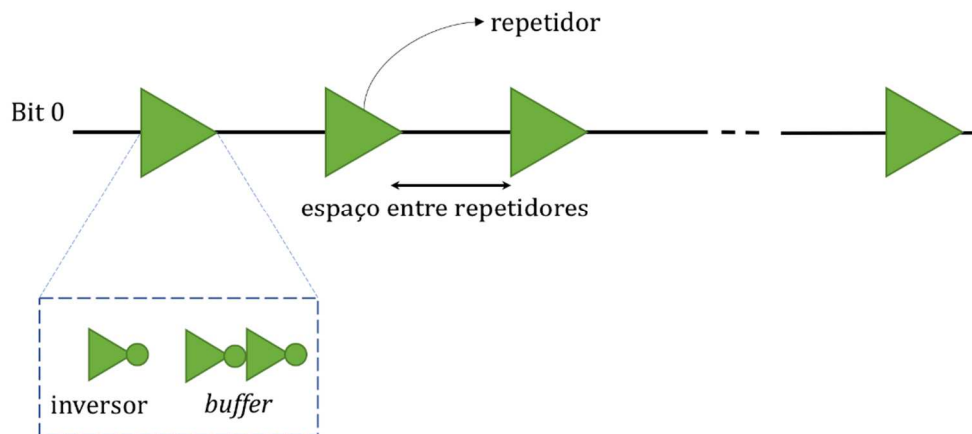


Figura 2-9 - Interconexão com repetidores

Na inserção do repetidor, geralmente a redução do atraso de interconexão é parcialmente compensada pelo atraso adicional dos repetidores inseridos. O atraso global do fio pode ser minimizado pela seleção de tamanhos de repetidores ótimos e espaçamento entre repetidores. Essa técnica é comumente empregada em processadores modernos. Quando adicionados de forma a otimizar o atraso, os repetidores tornam o atraso total do fio igual à média geométrica do atraso total do fio e do atraso do estágio do repetidor individual [34].

2.3.4 - Modelos de interconexão

Nesta seção serão apresentados os modelos equivalentes dos circuitos de interconexão de cobre e CNT que serão os materiais utilizados nos enlaces das NoCs estudadas nesse trabalho.

As dimensões associadas a uma seção transversal de interconexão são mostradas na Figura 2-10, onde os fios são representados por retângulos, L é o comprimento, W é a largura, T é a espessura, S é distância entre os condutores em uma mesma camada e H é a distância de separação entre camadas (espessura do isolante).

Entre camadas de metal, a largura, espessura e o espaçamento entre os fios podem ser alterados. O aumento da largura e do espaçamento entre fios, resulta em um menor atraso de propagação da interconexão, com a redução do produto RC. Porém, com o aumento da largura e do espaço das interconexões, a quantidade de fios por área reduz e, portanto, a largura de banda global do sistema também é reduzida. Assim, o projetista deve encontrar um equilíbrio de todos esses parâmetros para conseguir um menor atraso e uma maior largura de banda para o sistema [34].

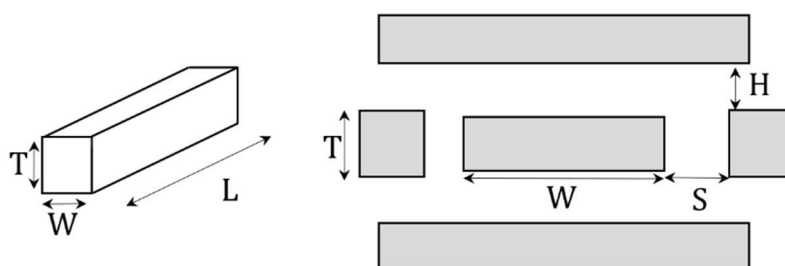


Figura 2-10 - Dimensões da interconexão

2.3.4.1. Modelo de Interconexão de cobre

Existem muitas opções para modelar o comportamento de uma conexão. O modelo mais simples é conhecido por modelo *lumped* que reúne todas as resistências em uma única equivalente e similarmente combina a capacitância total em um único capacitor equivalente.

O modelo *lumped* é impreciso e gera resultados pessimistas para longas interconexões. Assim, o modelo que melhor representa uma interconexão é o modelo π ou T, conhecido como modelo RC distribuído. Neste modelo, a interconexão é subdividida em segmentos. A precisão do modelo é determinada pelo número de segmentos N . A Figura 2-11 apresenta o modelo π_3 para a interconexão de cobre que será utilizado nesse trabalho. Esse modelo fornece um erro menor do que 3% [49].

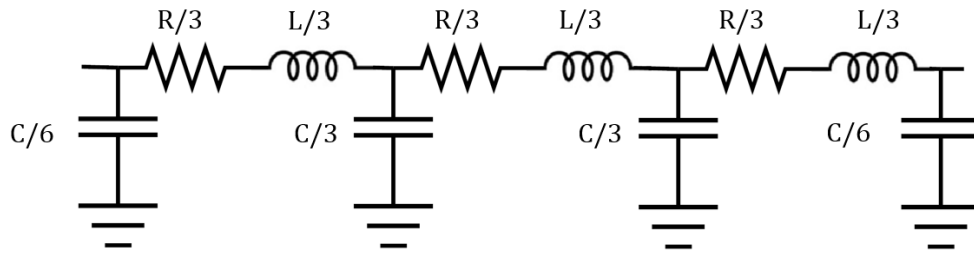


Figura 2-11 - Modelo de interconexão de cobre.

A seguir serão apresentados os parâmetros que constituem o modelo de interconexão de cobre apresentado.

2.3.4.1.1. Resistência do cobre

A resistência do cobre pode ser calculada pela Equação (2.4), onde ρ_{Cu} é a sua resistividade.

$$R_{Cu} = \frac{\rho_{Cu} \cdot L}{T \cdot W} \quad (2.4)$$

Em escala nanométrica, a resistividade do cobre é influenciada pela ocorrência dos fenômenos de espalhamento superficial e de espalhamento de contorno. O modelo proposto por Fuchs e Sondheimer (ρ_{FS}) e a teoria proposta por Mayadas e Shatkes (ρ_{MS}) quantificam esses fenômenos. Os parâmetros ρ_{FS} e ρ_{MS} são calculados pelas Equações (2.5) e (2.6), respectivamente, e o coeficiente α é dado pela Equação (2.7), onde ρ_o é a resistividade do cobre sem considerar os fenômenos citados anteriormente (*bulk*), l_o é o caminho médio livre dos elétrons do material de cobre, p_F é o parâmetro de espalhamento de Fuchs, D é o tamanho médio da região de depleção do contorno de grão, R é o coeficiente de reflexão no contorno com valores entre 0 e 1 e W é a largura do fio [9][44].

$$\rho_{FS} = \rho_o \left(1 + \frac{3}{4} (1 - p_F) \frac{l_o}{W} \right) \quad (2.5)$$

$$\rho_{MS} = \rho_o/3 \left[\frac{1}{3} - \frac{\alpha}{2} + \alpha^2 + \alpha^3 \cdot \ln \left(1 + \frac{1}{\alpha} \right) \right] \quad (2.6)$$

$$\alpha = \frac{l_o}{D} \cdot \frac{R}{1 - R} \quad (2.7)$$

2.3.4.1.2. Capacitância do cobre

A modelagem da capacitância do fio não é uma tarefa trivial. Porém, para fins de estimativa, algumas técnicas simples são aplicáveis e podem ser usadas. Dessa forma, conforme mostra a Figura 2-12, a capacitância do fio de cobre por unidade de comprimento pode ser modelada por quatro capacitores de placas paralelas, um para cada lado, e pela capacitância de borda (*fringing*). Os três componentes principais da capacitância total do fio (C_T), mostrados na Figura 2-12, estão relacionados pela Equação (2.8), onde C_a é capacitância de borda, C_b é a capacitância de placa paralela devido às camadas de metal superior e inferior e é proporcional a largura do fio, e C_c é a capacitância de acoplamento entre interconexões vizinhas no mesmo plano e é inversamente proporcional ao espaçamento de interconexão S [34].

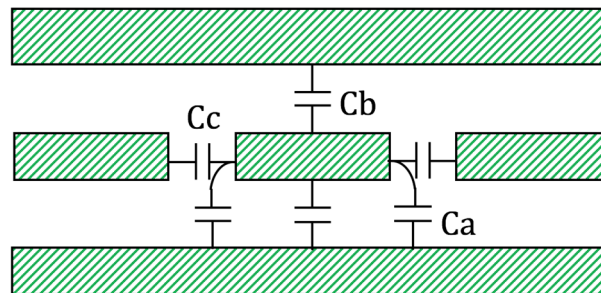


Figura 2-12 - Capacitâncias da interconexão.

$$C_T = C_a + 2 \cdot C_b \cdot W + \frac{C_c}{S} \quad (2.8)$$

A capacitância de placa paralela pode ser calculada pela Equação (2.9) [34], onde ϵ é a permissividade relativa para uma dada constante dielétrica.

$$C_b = \epsilon \frac{L}{H} \quad (2.9)$$

A capacitância de borda e a capacitância de acoplamento são mais difíceis de

serem calculadas e requerem solução via método dos elementos finitos para resultados mais precisos. Entretanto, para fins de estimativa e modelagem, as Equações (2.10) e (2.11) são computacionalmente eficientes e relativamente precisas [34].

$$C_a = \varepsilon \cdot L \left[\left(\frac{W}{H} \right) + 0,77 + 1,06 \left(\frac{W}{H} \right)^{0,25} + 1,06 \left(\frac{T}{H} \right)^{0,5} \right] \quad (2.10)$$

$$C_c = \varepsilon \cdot L \cdot S \left[0,03 \left(\frac{W}{H} \right) + 0,83 \left(\frac{T}{H} \right) - 0,07 \left(\frac{T}{H} \right)^{0,222} \right] \left(\frac{H}{S} \right)^{\frac{4}{3}} \quad (2.11)$$

2.3.4.1.3. Indutância do cobre

A indutância própria (L) e a mútua (M_{Cu}) da interconexão de cobre, em escala nanométrica, podem ser calculadas pelas Equações (2.12) e (2.13) [8] [9], onde μ_0 é a permeabilidade magnética do vácuo dada por $\mu_0 = 4\pi \cdot 10^{-7}$. Assim, a indutância total do cobre pode ser calculada pela soma das indutâncias própria e mútua.

$$L_{Cu} = \frac{\mu_0 \cdot l}{2\pi} \left[\ln \left(\frac{2l}{w+t} \right) + \frac{1}{2} + \frac{0,22(w+t)}{l} \right] \quad (2.12)$$

$$M_{Cu} = \frac{\mu_0 \cdot l}{2\pi} \left[\ln \left(\frac{2l}{s} \right) - 1 + \frac{s}{l} \right] \quad (2.13)$$

2.3.4.2. Modelo de Interconexão do SWCNT

Os CNTs são formados por uma lâmina de grafeno enrolada, denominada SWCNT (*Single-Walled Carbon Nanotube*), ou por um conjunto de lâminas concêntricas formando uma multicamada, denominada MWCNT (*Multi-Walled Carbon Nanotube*) [7]. A Figura 2-13 mostra a estrutura básica de um CNT.

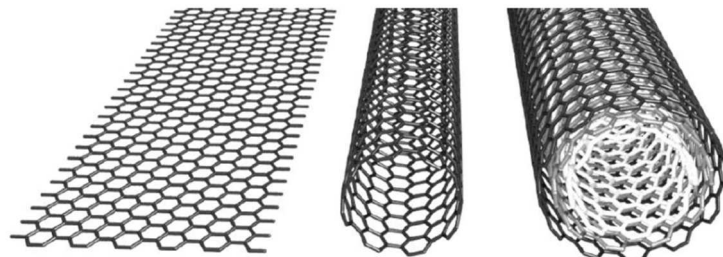


Figura 2-13 - Estrutura básica de um CNT. Lâmina de grafeno (esquerda), SWCNT (meio) e MWCNT (direita)[7].

O circuito equivalente do modelo de interconexão do SWCNT isolado é mostrado na Figura 2-14.

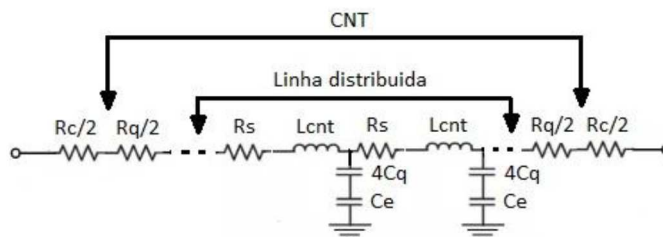


Figura 2-14 - Modelo de interconexão de SWCNT [50].

No circuito equivalente apresentado na Figura 2-14, R_C é a resistência de contato entre o metal e o nanotubo de carbono, R_q é a resistência quântica, R_S é a resistência de espalhamento, L_{CNT} é a indutância total do SWCNT, C_q é a capacitância quântica e C_e é a capacitância eletrostática [51][52]. A seguir serão apresentados os cálculos dos parâmetros dessas variáveis.

2.3.4.2.1. Resistência do SWCNT isolado

A resistência do SWCNT isolado possui três componentes principais: resistência de contato (R_C), resistência quântica (R_q) e a resistência de espalhamento (R_S).

As resistências de contato e quântica são fixas e não dependem do comprimento do nanotubo [53] [54]. Neste trabalho, a resistência de contato será considerada igual a 120 k Ω [57]. A resistência quântica pode ser calculada pela Equação (2.14), sendo dividida igualmente no modelo equivalente em cada lado dos contatos metal-nanotubo [6][7][55][56], onde h é a constante de Planck, e é a carga do elétron. A resistência de espalhamento depende do comprimento do nanotubo e pode ser calculada pela Equação (2.15), onde l_{CNT} é o comprimento do nanotubo de carbono e λ_{CNT} é o comprimento do caminho médio livre do SWCNT, que é tipicamente 1 μm .

$$R_q = \frac{h}{4e^2} = 6,45 \text{ k}\Omega \quad (2.14)$$

$$R_S = R_q \left(\frac{l_{CNT}}{\lambda_{CNT}} \right) \quad (2.15)$$

Para comprimentos de interconexão de SWCNT menores do que o caminho

médio livre ($l_{CNT} \leq \lambda_{CNT}$), o transporte de elétrons é exclusivamente balístico e a resistência independe do comprimento do nanotubo [7][57]. Porém, quando ($l_{CNT} \geq \lambda_{CNT}$), a resistência de espalhamento deve ser adicionada [7][8][9]. Assim, a resistência total do SWCNT isolado pode ser calculada Equação (2.16).

$$R_{CNT} = \begin{cases} R_C + R_q; & \text{se } l_{CNT} \leq \lambda_{CNT} \\ R_C + R_q + R_S; & \text{se } l_{CNT} \geq \lambda_{CNT} \end{cases} \quad (2.16)$$

2.3.4.2.2. Indutância do SWCNT isolado

O SWCNT possui duas componentes de indutância: indutância magnética (L_M) e indutância cinética (L_K). A indutância magnética é devido ao total de energia magnética resultante da corrente que flui no fio e pode ser expressa pela Equação (2.17). Para o cálculo dessa indutância, considera-se que o nanotubo é um fio muito fino, com diâmetro d_{CNT} e está situado a uma distância y do plano ligado ao terra. Já a indutância cinética surge da energia cinética armazenada em cada um dos canais condutores do CNT e pode ser calcula pela Equação (2.18), onde v_F é a velocidade de Fermi, cujo valor é igual a $8 \cdot 10^5$ m/s. Cada nanotubo possui quatro canais condutores em paralelo que não interagem entre si, o que resulta na indutância cinética efetiva dada por $L_K/4$ [7][8][9].

$$L_M = \frac{\mu}{2\pi} \left(\ln \frac{y}{d_{CNT}} \right) \quad (2.17)$$

$$L_K = \frac{\mu}{2e^2 v_F} \quad (2.18)$$

2.3.4.2.3. Capacitância do SWCNT isolado

O SWCNT possui duas capacitâncias de origens distintas: a capacitância eletrostática (C_E) e a capacitância quântica (C_Q). A capacitância eletrostática é ocasionada pelo carregamento de cargas do ambiente que envolve o fio, ou seja, os fios vizinhos e o plano terra, e pode ser calcula pela Equação (2.19). A capacitância quântica se refere a energia quântica armazenada no nanotubo quando este transporta corrente e pode ser calcula pela Equação (2.20).

$$C_E = \frac{2\pi\epsilon}{\ln \left(\frac{y}{d_{CNT}} \right)} \quad (2.19)$$

$$C_Q = \frac{2e^2}{h\nu_F} \quad (2.20)$$

Tendo em vista que o CNT é constituído por quatro canais condutores, a capacitância total do SWCNT isolado é dada pela Equação (2.21).

$$C_{CNT} = \frac{C_E \cdot 4C_Q}{C_E + 4C_Q} \quad (2.21)$$

2.3.4.3. Modelo de Interconexão do BCNT

O SWCNT *bundle*, ou BCNT (*Single-Walled Carbon Nanotube Bundle*) é formado por um conjunto de lâminas de grafeno enroladas individualmente e empacotadas em paralelo. Neste trabalho assume-se que todos os SWCNTs são idênticos, metálicos e que cada um possui o mesmo potencial [6] [58]. Dado que d_{CNT} é o diâmetro do nanotubo de carbono, no valor de 1 nm, e x a distância entre os centros de nanotubos adjacentes, o BCNT pode ser empacotado de forma densa, caso $x=d$, ou esparsa, caso $x > d$ [7][8][9][57]. Devido à força de Van der Waals, os nanotubos são separados por uma distância δ_{min} , que é de pelo menos 0,32 nm [7].

O número de nanotubos de carbono (n_{CNT}) disponível em um BCNT depende da largura e da espessura da interconexão. Assim, sabendo que n_W corresponde ao número de CNTs ao longo da largura do BCNT que pode ser obtido pela Equação (2.22), e que n_T corresponde ao número de CNTs ao longo da espessura do BCNT que pode ser obtido pela Equação (2.23), tem-se que o n_{CNT} pode ser calculado pela Equação (2.24).

$$n_W = \left\lfloor \frac{W - d_{CNT}}{x} \right\rfloor \quad (2.22)$$

$$n_T = \left\lfloor \frac{T - d_{CNT}}{(\sqrt{3}/2)x} \right\rfloor + 1 \quad (2.23)$$

$$n_{CNT} = \begin{cases} n_W n_T - \frac{n_T}{2}; & \text{se } n_T \text{ par} \\ n_W n_T - \frac{n_T - 1}{2}; & \text{se } n_T \text{ ímpar} \end{cases} \quad (2.24)$$

Após a definição desses parâmetros iniciais, pode-se finalmente calcular os valores de resistência, indutância e capacitância para esse modelo.

2.3.4.3.1. Resistência do BCNT

O BCNT possui resistência equivalente menor do que o SWCNT isolado, o que permite alcançar desempenho semelhante ao das interconexões de cobre [7][8][9]. Assim, a resistência total do BCNT pode ser calculada pela Equação (2.25).

$$R_{bundle} = \frac{R_{CNT}}{n_{CNT}} \quad (2.25)$$

Assim como para o SWCNT isolado, neste trabalho considera-se a resistência de contato para o BCNT como sendo ideal.

2.3.4.3.2. Indutância do BCNT

A indutância do BCNT é dada pela combinação em paralelo de cada indutância do SWCNT que forma o BCNT e pode ser calculada pela Equação (2.26).

$$L_{bundle} = \frac{L_{CNT}}{n_{CNT}} \quad (2.26)$$

2.3.4.3.3. Capacitância do BCNT

A capacitância do BCNT é composta por duas componentes: capacitância quântica (C_Q^{bundle}) e capacitância eletrostática (C_E^{bundle}). Essas capacitâncias podem ser calculadas pelas Equações (2.29) e (2.30) respectivamente, onde C_{En} , calculado pela Equação (2.27), representa a capacitância entre placas paralelas próximas de SWCNTs e C_{Ef} , calculada pela Equação (2.28) representa a capacitância entre placas paralelas afastadas de SWCNTs. Por fim, a capacitância total do BCNT (C_{bundle}) é expressa pela Equação (2.31).

$$C_{En} = \frac{2\pi\epsilon}{\ln\left(\frac{S}{d_{CNT}}\right)} \quad (2.27)$$

$$C_{Ef} = \frac{2\pi\epsilon}{\ln\left(\frac{S+W}{d_{CNT}}\right)} \quad (2.28)$$

$$C_Q^{bundle} = C_Q^{SWCNT} \cdot n_{CNT} \quad (2.29)$$

$$C_E^{bundle} = 2C_{En} + \frac{n_W - 2}{2} C_{Ef} + \frac{3(n_T - 2)}{5} C_{En} \quad (2.30)$$

$$C_{bundle} = \frac{C_E^{bundle} \cdot C_Q^{bundle}}{C_E^{bundle} + C_Q^{bundle}} \quad (2.31)$$

Conforme pode ser observado na Equação (2.31), para valores grandes de n_{CNT} , o efeito da capacitância quântica é pequeno. Com isso, a capacitância total do BCNT pode ser aproximada pela sua capacitância eletrostática [6][7][44][52].

Segundo Srivastava et. al. [7], os nanotubos no interior do BCNT são blindados eletrostaticamente dos condutores de terra, podendo ser desprezados no cálculo da capacitância eletrostática. Assim, neste cenário, apenas os CNTs de borda são levados em consideração. No entanto, a Equação (2.30) não reproduz essa realidade, pois considera a capacitância eletroestática de todos CNTs e não somente os de borda. Assim, neste trabalho foi adotado que o valor da capacitância total do BCNT por unidade de comprimento é igual ao do cobre, considerando a mesma seção transversal do fio, com base nas observações e análises realizadas por Pasricha *et al.* [59].

3 - ESTIMATIVA DO CONSUMO DE ENERGIA DE REDES-EM-CHIP

3.1 - CONSUMO DE ENERGIA DA INTERCONEXÃO

Conforme visto anteriormente, a evolução dos projetos SoC resultou na arquitetura de comunicação baseada em NoCs. A adoção desse tipo de arquitetura tem por objetivo alcançar melhores métricas em um dispositivo tais como, redução do consumo de energia, maior largura de banda e menor latência. Por outro lado, a redução das dimensões dos dispositivos e o estreitamento das interconexões eleva o calor dissipado do sistema, devido à redução do espaço entre fios e ao aumento da frequência de operação do sistema. Nesta seção, é apresentado o modelo adotado para obter a energia consumida pela interconexão que conecta dois roteadores em uma NoC.

3.1.1 - Obtenção do consumo de energia da interconexão

Considere o modelo de um único fio de interconexão mostrado na Figura 3-1.

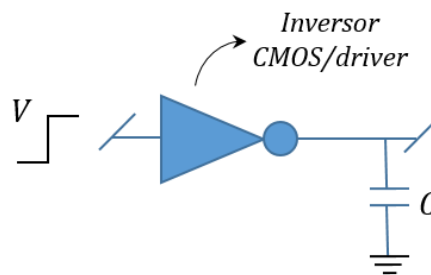


Figura 3-1 - Modelo simples de interconexão.

O consumo de energia da interconexão é comumente definido pela Equação (3.1). Esse modelo é baseado no pressuposto de que a energia é consumida em cada transição de aumento ou queda de um evento, onde C é a capacitância total da carga, dada pela soma da capacitância do fio e a capacitância da porta de entrada do inversor CMOS que conduz o fio (*driver*), e V é a tensão de alimentação da fonte [60].

$$E = \frac{1}{2} CV^2 \quad (3.1)$$

Entretanto, a energia dissipada ocasionada pelos efeitos da capacitância de acoplamento entre fios e da capacitância entre placas paralelas, não são considerados neste modelo. Esses efeitos são observados de forma proeminente quando os fios estão comutando para diferentes valores de saída ou quando um dos fios comuta de nível de tensão, enquanto o outro não comuta, resultando em diferentes valores de tensão de saída. Para tecnologias abaixo de 90nm, os valores típicos das distâncias que separam

os fios são reduzidos e os valores da capacitância de acoplamento são cada vez maiores. Assim, o efeito de acoplamento não pode ser desconsiderado se a capacitância de acoplamento for próxima ou maior que a capacitância da carga do fio [10].

Conforme visto anteriormente, a modelagem da capacitância do fio não é uma tarefa trivial. Assim, neste trabalho a energia consumida pela interconexão para transmitir um *bit* (E_{lbit}) foi obtida por meio do *software* LTspice [61], utilizando a capacitância total do modelo de interconexão dada pela Equação (2.9). Ainda, por questões de simplicidade, presume-se que a energia de uma interconexão composta por N fios pode ser calculada pela Equação (3.2). Ou seja, E_{link} é a energia consumida pela interconexão da NoC, para transmitir um *flit* e

$$E_{link} = NE_{lbit} \quad (3.2)$$

3.2 - CONSUMO DE ENERGIA DO ROTEADOR

3.2.1 - Arquitetura do Roteador Nanoelétrônico

Neste trabalho, no intuito de verificar o consumo de energia em NoCs baseadas em dispositivos nanoelétrônicos, será utilizado o roteador completamente baseado na tecnologia SET [16]. O esquemático completo da arquitetura do roteador nanoelétrônico é mostrado na Figura 3-2.

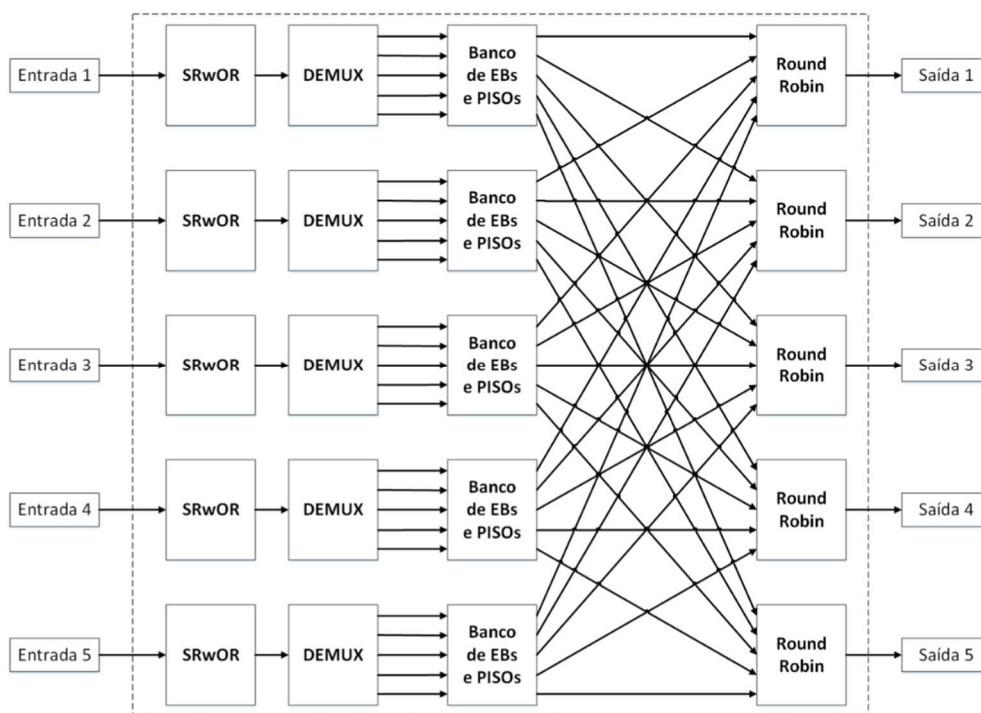


Figura 3-2 - Esquemático completo do roteador nanoelétrônico [16].

O roteador em questão foi desenvolvido para topologia em malha, com cinco entradas e cinco saídas, composto por registradores, demultiplexadores (DEMUX), *buffers*, conversores e componentes de arbitragem. Cada entrada do roteador possui um bloco registrador de deslocamento com registrador de saída (SRwOR - *Shift-Register with Output Register*) do tipo SIPO (*Serial-in, Parallel-out*), para realizar a conversão serial-paralelo de dados, onde um *flit* de 8 bits é processado e o endereçamento é extraído. Em seguida, o endereçamento é utilizado pelo bloco DEMUX que seleciona a saída correta e armazena o *flit* no buffer elástico (EB - *Elastic Buffer*).

O bloco denominado por “Banco de EBs e PISOs” possui cinco EBs e cinco registradores de deslocamento do tipo que converte uma entrada paralela em uma saída serial (PISO - *Parallel-in, Serial-out Register*). Os EBs realizam o armazenamento temporário dos dados e executam o protocolo de controle de fluxo *handshake*. Após o pacote ser armazenado no EB, um sinal é enviado ao árbitro solicitando o uso do recurso de saída. Após a confirmação da solicitação pelo árbitro que utiliza o protocolo de arbitragem *Round and Robin* (RoR), o *flit* é novamente serializado por meio do registrador PISO e finalmente enviado ao canal de saída correspondente. A Figura 3-3 apresenta o diagrama de blocos simplificado para transferir um *flit* da entrada até a saída do roteador nanoeletrônico.

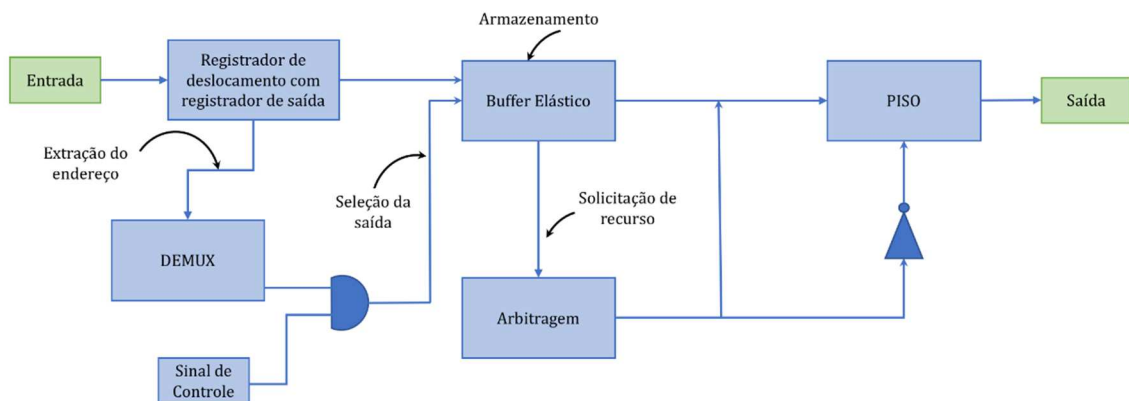


Figura 3-3 - Diagrama de blocos simplificado do fluxo de um flit.

O roteador em questão foi desenvolvido para o processamento de um *bit* a cada 9 ciclos de *clock*, em uma rede sem contenção de dados. Esse tempo de processamento inclui a paralelização e posterior serialização do *bit*, execução do protocolo de controle de fluxo e arbitragem.

3.2.2 - Obtenção do consumo de energia do roteador

O roteador nanoeletrônico [16] foi desenvolvido apenas com portas NAND nanoeletrônicas. Assim, a potência e a área dissipada pelo circuito do roteador foram

calculadas a partir da área e potência de uma porta NAND nanoeletrônica (nanoNAND). Uma nanoNAND possui 172 nm² de área e dissipa 220 pW de potência considerando uma tensão de operação de 0,9V e frequência de 1 GHz [16]. A Tabela 3.1 apresenta a quantidade total de nanoNANDs em cada módulo, a área e a potência dissipada para cada módulo, bem como para o roteador completo [16].

Tabela 3.1 – Área e potência dos módulos do roteador nanoeletrônico [16].

	Total de nanoNAND	Área	Potência Dissipada
nanoNAND	1	172 nm ²	220 pW
SRwOR	368	63296 nm ²	80960 pW
DEMUX	53	9116 nm ²	11660 pW
EB	2318	398696 nm ²	509960 pW
PISO	209	35948 nm ²	45980 pW
RoR	193	33196 nm ²	42460 pW
Roteador 5x5	69619	12 μm ²	~15 μW

Para fins de comparação, caso a mesma arquitetura de roteador fosse construída com transistores CMOS com canal de 22nm, a potência dissipada pelo roteador seria de 32,72 mW [16]. Dessa forma, como pode ser observado, os dispositivos baseados em tecnologia nanoeletrônica são promissores para redução do consumo de energia dos circuitos.

Para calcular a energia total consumida pelo roteador para transmitir um *bit*, a energia de cada componente que compõe o roteador é somada, conforme mostrado na Equação (3.3). Para o roteador nanoeletrônico em questão, a energia total consumida para transferir um *bit* da entrada até a saída do roteador nanoeletrônico é de aproximadamente 26 fJ, enquanto que para o roteador CMOS com a mesma arquitetura do roteador nanoeletrônico, a energia consumida é de aproximadamente 56 pJ.

$$E_{roteador} = E_{SRWOR} + E_{DEMUX} + 5 \cdot E_{EB} + E_{RoR} + 5 \cdot E_{PISO} \quad (3.3)$$

3.3 - GERAÇÃO DE TRÁFEGO

A geração de tráfego define a forma de transmissão de dados em uma rede e, com isso, os nós de origem e destino são determinados, bem como o fluxo de mensagens entre eles, descrevendo a comunicação e o funcionamento de uma dada aplicação. A definição da geração de tráfego e do algoritmo de roteamento permite ao projetista a avaliação do desempenho da rede, por meio da validação dos requisitos da aplicação.

Para geração espacial de tráfego, alguns padrões são utilizados para estressar a arquitetura de comunicação e o seu algoritmo de roteamento. Nesta seção,

inicialmente, será apresentado o padrão de tráfego uniforme que não possui um destino específico. Em seguida, serão apresentados alguns padrões que apresentam um destino específico para cada núcleo gerador de tráfego. Estes tipos de padrões podem apresentar maior localidade temporal, ou seja, uma maior afinidade de comunicação entre núcleos da rede. Vale ressaltar que para uma aplicação ter uma melhor localidade temporal, não é necessário que ela tenha uma melhor localidade espacial [31].

3.3.1 - Tráfego uniforme aleatório

No tráfego uniforme aleatório, cada fonte possui probabilidade igual de enviar pacotes para cada destino. Esse padrão de tráfego é comumente utilizado para a avaliação de uma rede, pois é simples de implementar, não faz pressupostos sobre a aplicação e é analiticamente tratável. Como os nós de origem não diferenciam os nós de destino que estão mais próximos daqueles que estão mais longe, o tráfego aleatório uniforme não explora a localidade de comunicação e, portanto, pode ser utilizado como um estudo de caso, em condições não ideais.

3.3.2 - Tráfego permutação de bit

No tráfego permutação de bit, cada fonte envia todo o seu tráfego para um único destino. Como esse tipo de tráfego concentra carga em pares individuais fonte-destino, eles tendem a estressar o equilíbrio de carga de uma topologia e seu algoritmo de roteamento. As permutações de bits são uma subclasse de permutações nas quais o endereço de destino é calculado permutando os bits do endereço de origem. Detalhes sobre como gerar esses padrões de tráfego são mostrados em [24]. Neste trabalho, serão utilizados os tráfegos de permutação de *bits* do tipo *Bit Complement* e *Bit Rotation*.

3.3.3 - Tráfego *Nearest Neighbor*

O tráfego *Nearest Neighbor* é comumente usado para avaliar o impacto da localidade de comunicação sobre o desempenho e consumo de energia da rede no *chip* [62]. Uma porcentagem fixa de tráfego vai para os vizinhos mais próximos que se encontram a uma distância de raio r e o resto do tráfego segue uma distribuição uniforme e aleatória.

3.4 - MODELANDO A LOCALIDADE ESPACIAL DE COMUNICAÇÃO NoC USANDO A REGRA DE RENT

A regra de Rent é um padrão experimental observado em projetos VLSI que descreve a estrutura de comunicação entre portas lógicas em um *chip*, na qual as interconexões são localizadas de forma a minimizar a potência e latência do fio [63].

Usando derivações baseadas na regra de Rent, a distribuição de comprimento de fio (WLD - *Wire Length Distribution*) de um circuito pode ser estimada a partir do expoente e coeficiente de Rent, p e k respectivamente [64]. Essa distribuição descreve a localidade de comunicação no circuito e, portanto, é fundamental para projetos VLSI, pois está relacionada a muitas propriedades do sistema, como a área do *chip*, atraso de sinal, consumo de energia e o roteamento dos fios [65] [17].

Em SoCs, informações semelhantes são fornecidas pela Distribuição de Probabilidade de Comunicação (CPD - *Communication Probability Distribution*) das aplicações. A CPD descreve a probabilidade dos pacotes percorrerem uma determinada distância, em uma NoC, para um determinado padrão de tráfego. Esta distribuição espacial está diretamente relacionada com o consumo de energia de uma aplicação, porque quanto maior a distância percorrida pelos pacotes, maior é a energia consumida. Tendo em vista que uma NoC pode chegar a consumir de 30 a 40% do orçamento de potência de um *chip* [66][67], é desejável que a distância percorrida pelos pacotes seja a menor possível, para reduzir o consumo energético do *chip*.

Neste trabalho, a CPD será utilizada para estudar a distribuição espacial do tráfego e o consumo de energia de uma NoC. Dado que o padrão de comunicação de muitas aplicações paralelas segue a regra de Rent [68], um gerador de padrões de tráfego foi proposto [17], onde a probabilidade de comunicação entre processadores é derivada diretamente da regra de Rent, produzindo CPDs com alta localidade de tráfego. Esse gerador foi desenvolvido com o objetivo de simular o tráfego das aplicações em uma NoC de maneira rápida e simples. Ainda, a partir dessas CPDs, um modelo para prever o consumo de energia em uma NoC também foi proposto [17]. Este modelo será utilizado posteriormente nesta dissertação, para realizar a comparação do consumo de energia de diferentes NoCs. Essa abordagem foi escolhida, pois não requer simulação e poderia ser usada nas fases iniciais do projeto de NoC, podendo inclusive auxiliar no projeto de aplicações eficientes em energia e em melhores técnicas de mapeamento de aplicação [17].

3.4.1 - Geração de tráfego utilizando a regra de Rent

O gerador sintético proposto [17] segue a regra de Rent e emprega o tráfego sintético dessa regra como um modelo genérico de comunicação em aplicações paralelas. Esse gerador tem por objetivo avaliar NoCs com cargas de trabalho que reproduzem as propriedades espaciais do tráfego real.

Em VLSI, a probabilidade de um fio conectar dois terminais separados por uma distância de Manhattan d é dada pela Equação (3.4) [1][17].

$$P(d) = \frac{(1 + d(d - 1))^p - (d(d - 1))^p + (d(d + 1))^p - (1 + d(d + 1))^p}{4d} \quad (3.4)$$

A equação acima é utilizada para definir a probabilidade de comunicação entre dois núcleos, onde d corresponde ao número de saltos no caminho mais curto entre a fonte e o destino. Assim, o tráfego de cada nó de origem pode ser gerado a partir da probabilidade dada pela Equação (3.4) para cada nó de destino possível. A repetição desse processo para todos os nós da rede resulta no tráfego que segue regra de Rent.

Uma propriedade interessante deste método é a capacidade de gerar vários padrões de tráfego variando apenas o expoente Rent. Assim, como o expoente do Rent está relacionado à localidade de comunicação e à complexidade das aplicações, é possível analisar uma NoC em vários cenários de aplicações.

3.4.2 - Distribuição de probabilidade de comunicação da regra de Rent

A fórmula para calcular a CPD do tráfego sintético da regra de Rent é dada pela Equação (3.5), onde Γ é o coeficiente de normalização dado pela Equação(3.6) [17] e N corresponde ao número de nós na rede.

$$CPD(d) = \Gamma P(d) * \sum_{i=1}^{2\sqrt{N}-2} (\sqrt{N} - i)(\sqrt{N} + i - d), \quad (3.5)$$

$$para \ 0 < (\sqrt{N} + i - d) \leq \sqrt{N},$$

$$\Gamma = 1 / \sum_{i=1}^{2\sqrt{N}-2} CPD(d) \quad (3.6)$$

A Figura 3-4 mostra a CPD dos padrões de tráfego apresentados nesta seção, para uma rede 8x8 com topologia em malha. A Figura 3-4 (a) mostra a CPD para o tráfego uniforme aleatório. As CPDs do *Bit Complement* e *Bit Rotatiton* são mostradas nas Figura 3-4 (b) e (c), respectivamente. Estas distribuições são consideravelmente diferentes umas das outras, bem como do tráfego aleatório uniforme. A CPD do tráfego *Nearest Neighbor* com $r = 1$ e fator de localidade de 50% é mostrada na Figura 3-4 (d). Por fim, a CPD para o tráfego sintético proposto [17] é mostrada na Figura 3-4 (e).

Como pode ser observado, os padrões de tráfego sintético encontrados na literatura são úteis para analisar a rede, mas têm pouca ou nenhuma semelhança com

o tráfego real, pois as cargas de trabalho geradas por esse tipo de tráfego possuem uma localidade de comunicação pobre. Assim, o tráfego gerado pela regra de Rent, torna-se uma opção atraente, pois possui maior localidade de comunicação do que a maioria dos tráfegos sintéticos.

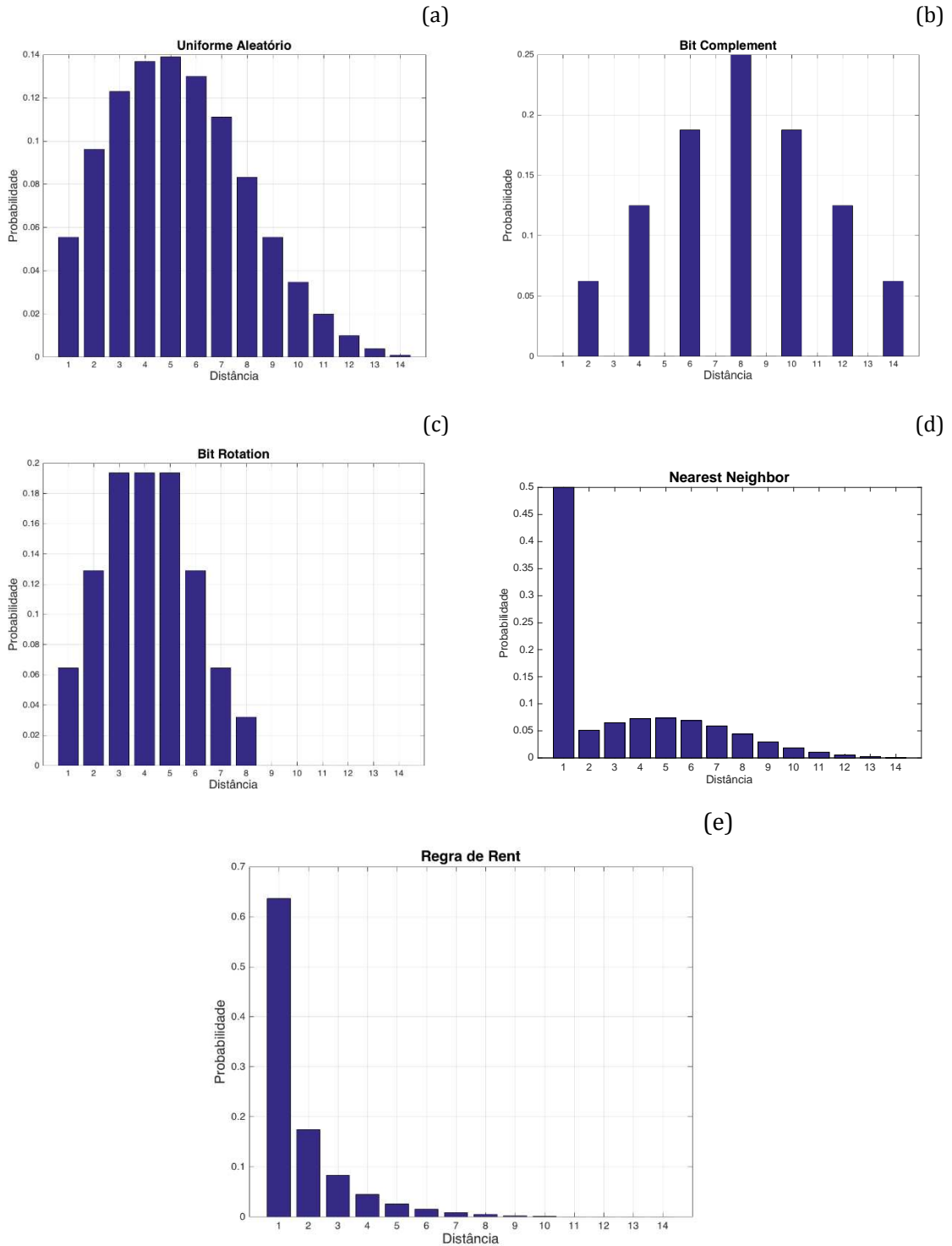


Figura 3-4 - CPD para diferentes tipos de padrão de tráfego em uma rede 8x8 com topologia em malha. (a) Uniforme aleatório. (b) Bit complement. (c) Bit rotation. (d) Nearest Neighbor com fator de localização de 50 %. (e) Regra de Rent com expoente de 0,75.

3.5 -MODELO ANALÍTICO PARA CÁLCULO DO CONSUMO DE ENERGIA EM NoCS

Analisar o consumo de energia de uma NoC usando simulações pode ser computacionalmente caro e proibitivo, especialmente em sistemas com cargas de trabalho orientadas a aplicativos ou em circuitos muito grandes. Nesta seção é apresentado o modelo analítico proposto por Bezerra [17] para prever o consumo de energia de uma NoC, baseado na sua CPD, que não requer simulações por computador. Este modelo é destinado a redes diretas em que o comprimento dos fios entre núcleos é o mesmo, ou seja, o tamanho do salto é igual, como por exemplo em redes com topologia *mesh* ou toróide. Porém, esse modelo pode ser facilmente estendido e adaptado para outras topologias [17].

A energia média consumida por um *flit* para percorrer um caminho de comprimento d em uma NoC é dada pela Equação (3.7), onde E_{link} e $E_{roteador}$ são a energia consumida pelo *flit* para atravessar um enlace e um roteador, respectivamente, e d corresponde ao número de saltos percorridos na transmissão do *flit* do nó de origem até o seu nó de destino [17].

$$E_{flit}(d) = d \cdot E_{link} + (d + 1) \cdot E_{roteador} \quad (3.7)$$

A energia total consumida por uma NoC, a qual representa a comunicação entre os núcleos de uma aplicação, é obtida pela Equação (3.8). Como pode ser observado nessa equação, a energia total consumida é dada pela soma da energia média de um *flit* sobre todas as distâncias de comunicação que ligam os nós dentro da topologia de rede escolhida, ponderadas pela probabilidade do pacote viajar essa distância. Este valor é então multiplicado pelo número de *flits* por pacote (N_{flits}) e o número total de pacotes ($N_{pacotes}$). Na Equação (3.8), assume-se que o número de *flits* por pacote é constante [17]. Neste trabalho, a constante E_{link} foi obtida por meio do software LTspice e das considerações descritas na seção 3.1, enquanto que a constante $E_{roteador}$ foi obtida por meio da adaptação do cálculo de potência da arquitetura do roteador [16].

$$E_{total} = N_{pacotes} \cdot N_{flits} \cdot \sum_{d=1}^{max} E_{flit}(d) \cdot CPD(d) \quad (3.8)$$

Dado o expoente de Rent, a CPD do tráfego de Rent pode ser obtida diretamente das Equações (3.4) e (3.5). Com essas informações, o consumo de energia da NoC de uma aplicação pode ser facilmente previsto a partir da Equação (3.8). A capacidade do modelo proposto [17] de calcular o consumo de energia para o tráfego Rent com base

em um único parâmetro de aplicação, pode simplificar e acelerar significativamente a análise de energia de uma NoC, o que apresenta ser uma vantagem sobre outras abordagens de cálculo do consumo de energia em uma NoC [70][71][72]. O gerador de tráfego baseado na regra de Rent, reproduz a CPD de padrões de tráfego para aplicações reais. Assim, esse método pode ser usado como uma maneira simples de avaliar uma NoC sob uma variedade de cenários, sem ter que recorrer a cargas de trabalho orientadas por aplicativos [17].

Uma limitação potencial deste método é a suposição de que a energia utilizada para a comunicação é proporcional à distância percorrida pelos pacotes. Isto é aproximadamente verdadeiro para a maioria das aplicações que utilizam NoCs e é comumente adotado na literatura para fins de simplificação [70][71][72]. Além disso, a contenção na rede pode aumentar a energia dinâmica e estática que não são contabilizadas pelo modelo [17].

4 - METODOLOGIA

4.1 - INTRODUÇÃO

Conforme visto anteriormente, uma NoC é composta por roteadores interligados por meio de enlaces. Dessa forma, para calcular o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos, a metodologia desenvolvida é dividida em três partes. A primeira parte refere-se à obtenção da energia das interconexões de cobre e BCNT, isoladamente, a segunda parte refere-se à obtenção da energia do roteador e, por fim, a terceira parte utiliza os dados de energia da interconexão e do roteador obtidos nas etapas anteriores, para calcular o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos.

4.2 - OBTENÇÃO DO CONSUMO DE ENERGIA DAS INTERCONEXÕES GLOBAIS DE COBRE E BCNT

A Figura 4-1 mostra o fluxograma para a obtenção do consumo de energia das interconexões, seja de cobre ou BCNT. O modelo distribuído π 3 foi escolhido para simular esses dois materiais de interconexão, por meio do *software* LTspice. Dado que a interconexão entre dois roteadores é do tipo global, primeiramente, faz-se necessário definir os parâmetros desse tipo de interconexão, com a relação ao nó de tecnologia utilizado. Assim, as tecnologias de 22 nm, 45 nm, 65 nm e 90 nm foram estudadas, em dois cenários diferentes. O primeiro cenário utiliza os parâmetros disponibilizados pela INTEL, na fabricação dos seus *chips*, enquanto que o segundo utiliza os parâmetros disponibilizados pelo ITRS (*International Technology Roadmap for Semiconductors*) [33].

Para o cenário da INTEL, a sexta camada da pilha de interconexões das tecnologias estudadas foi adotada como a camada responsável pela interligação dos componentes dentro de uma rede. Com a definição da camada, a largura e espessura da interconexão global foram obtidas para todos os nós de tecnologia, por meio de literatura científica previamente disponibilizada [73][74][75][76][77]. No cenário do ITRS, a largura e espessura da interconexão global, e a resistividade do material de cobre foram obtidas em função do nó de tecnologia [33]. Tanto no caso da INTEL, quanto no caso do ITRS, os valores da permissividade do material utilizado entre as camadas foram retirados dos relatórios do ITRS.

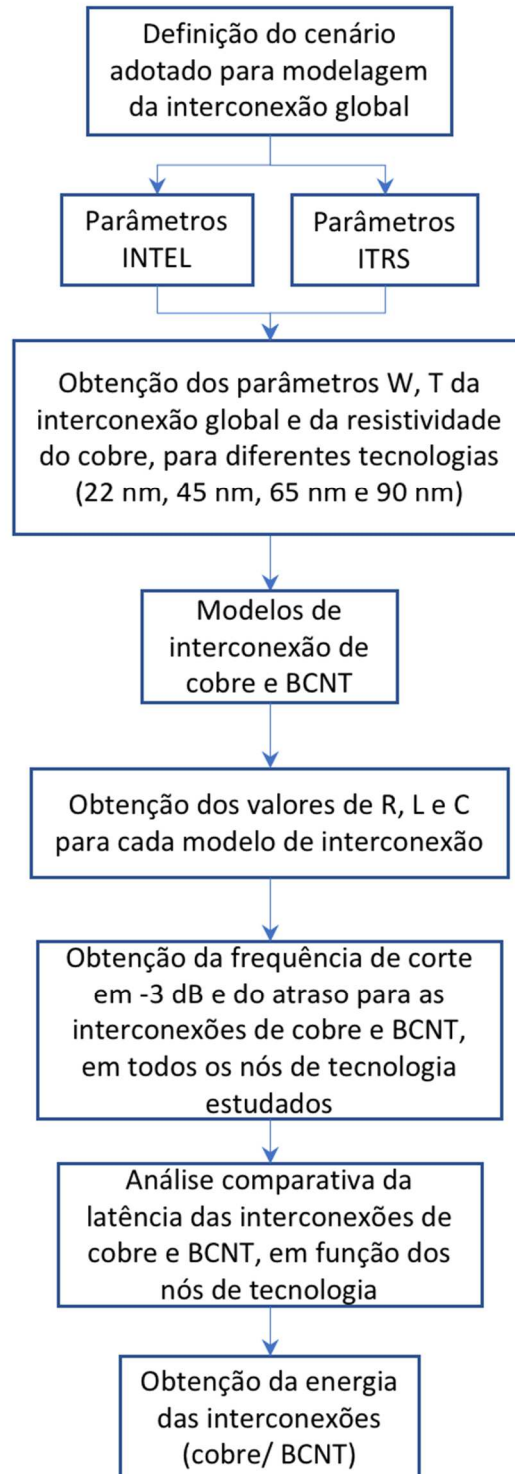


Figura 4-1 - Fluxograma das etapas para obtenção do consumo de energia para as interconexões de cobre e BCNT.

Após o levantamento da seção transversal do fio e da sua resistividade, a partir dos modelos de interconexão do cobre e do BCNT apresentados nas seções 2.3.4.1 e 2.3.4.3, respectivamente, os valores da resistência, indutância e capacitância foram

obtidos para cada um desses dois materiais. Para realizar esses cálculos, o comprimento escolhido para interconexão global foi de 1 mm. Assim, a partir da simulação do circuito da Figura 4-2, a frequência de corte em - 3 dB e a potência de cada interconexão foram encontradas, e com isso, o tempo de atraso para propagar um *bit* foi calculado. A tensão de entrada utilizada (V_{in}) foi uma onda quadrada de 0,9 V de amplitude, e o capacitor da carga (C_{carga}) utilizado possui 100 aF.

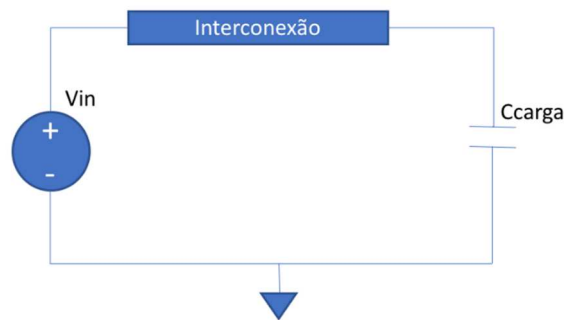


Figura 4-2 - Circuito utilizado na simulação das interconexões.

Em seguida, a análise comparativa da latência para as interconexões de cobre e BCNT foi realizada, em função de todas as tecnologias estudadas e, finalmente, o consumo de energia das interconexões globais para esses materiais foi obtido.

4.3 -OBTENÇÃO DO CONSUMO DE ENERGIA DO ROTEADOR

A Figura 4-3 mostra o fluxograma para a obtenção do consumo de energia do roteador baseados em dispositivos CMOS ou nanoeletrônicos.

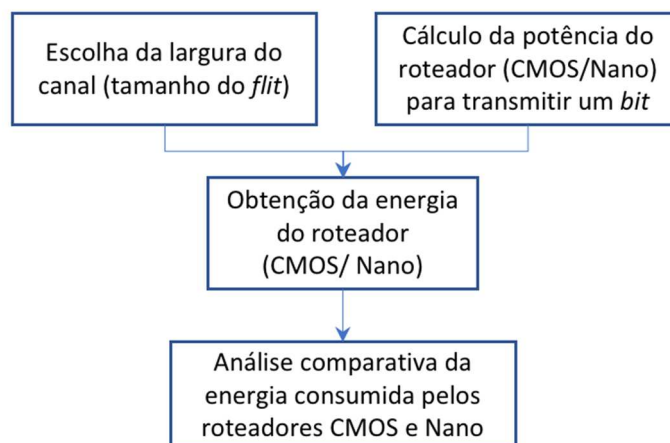


Figura 4-3 - Fluxograma das etapas para obtenção do consumo de energia do roteador baseado em dispositivos CMOS ou nanoeletrônicos.

Para obter o consumo de energia do roteador, primeiramente calculou-se a potência consumida para processar um *bit* da entrada até a saída do roteador

nanoeletrônico [16], utilizando a Equação (3.3). Para fins de comparação, a partir dos dados obtidos para o roteador nanoeletrônico [16], a potência dissipada para processamento de um *bit* foi calculada para o roteador CMOS.

Após obter a potência do roteador, com o objetivo de calcular o consumo de energia dos roteadores para diferentes tamanhos de *flit*, o consumo de potência obtido foi adaptado, para que fosse possível sua utilização neste trabalho. Assim, o consumo de energia dos roteadores nanoeletrônico e CMOS foram obtidos para diferentes tamanhos de *flit*, utilizando o modelo adaptado e considerando que o roteador necessita de três ciclos de *clock* para processamento de um *flit* e, com isso, foi possível comparar a energia consumida por esses roteadores.

4.4 -OBTENÇÃO DO CONSUMO DE ENERGIA DE UMA NOC

A Figura 4-4 mostra o fluxograma para a obtenção do consumo de energia de NoCs.

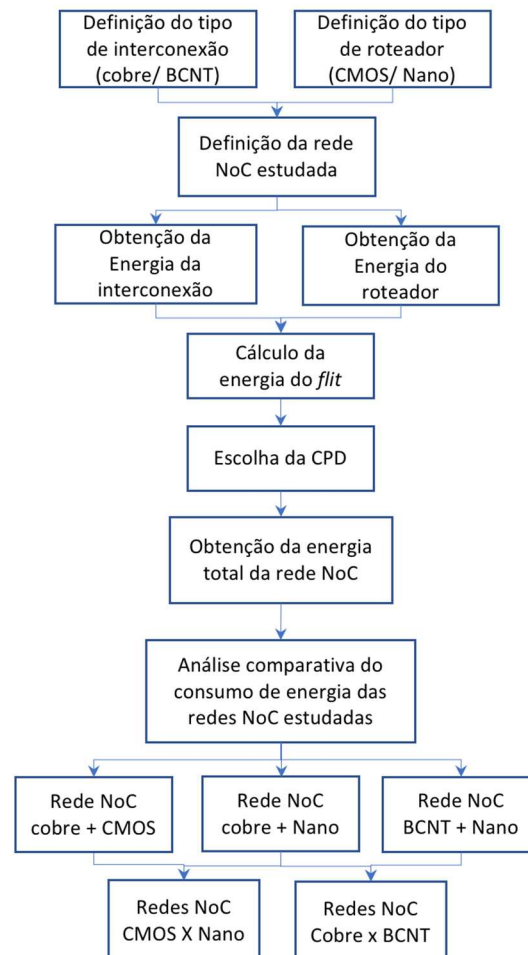


Figura 4-4 - Fluxograma das etapas para obtenção do consumo de energia de NoCs.

Nesse trabalho, o modelo analítico descrito na seção 3.5 foi adotado, no intuito de realizar o estudo sobre o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos. Dessa forma, com o objetivo de realizar esse estudo, outras redes foram construídas, a fim de comparar a dissipação de energia entre elas e, com isso, possibilitar a avaliação da NoC baseada em dispositivos nanoeletrônicos.

Tendo em vista que uma NoC é constituída basicamente por roteadores e interconexões, a simples troca do tipo de tecnologia utilizada por esses elementos pode ter grande impacto no consumo de energia da comunicação da rede. Assim, neste trabalho, as redes estudadas foram escolhidas com base na definição dos tipos de tecnologia utilizados nas interconexões e nos roteadores.

Para o estudo das redes escolhidas, a frequência de operação do *clock* adotada foi de 1 GHz e, para os cálculos analíticos do consumo de energia, adotou-se o modelo de rede sem contenção, onde o roteador necessita de três ciclos de *clock* para processar um *flit*, enquanto que a interconexão entre dois roteadores necessita de um ciclo de *clock* para transmitir o *flit*. Este caso é semelhante ao roteador SCC de 48 núcleos da Intel sem contenção [79].

Para calcular a energia consumida para transmitir um *flit* do seu nó de origem até o seu nó de destino, a Equação (3.7) foi utilizada. Entretanto, para utilizar essa equação, primeiramente foi necessário obter o consumo de energia consumido pelo *flit* para atravessar uma interconexão e um roteador. Esses dados foram obtidos seguindo as etapas da metodologia, apresentadas anteriormente neste capítulo.

Após encontrar o consumo de energia para transmitir o *flit* na NoC, a Equação (3.8) foi utilizada para obter o consumo de energia total da NoC. Como pode ser observado nessa equação, alguns parâmetros de rede precisam ser previamente definidos e encontrados. Assim, foi necessário definir a quantidade do número de pacotes e do número de *flits* por pacote utilizados pela rede. Com isso, a energia total consumida pela NoC foi obtida pela soma da energia média de um *flit* sobre todas as distâncias de comunicação que ligam os nós dentro da topologia de rede escolhida, ponderadas pela CPD utilizadas, a qual depende do padrão de tráfego escolhido.

Assim, o modelo analítico foi aplicado para avaliar o consumo de energia de uma NoC baseada em dispositivos nanoeletrônicos e, posteriormente, o de uma NoC baseada em dispositivos CMOS, ambas construídas com interconexões de cobre. Em seguida, foi realizada uma análise comparativa entre o consumo de energia das redes construídas com roteadores CMOS e das redes construídas com roteadores nanoeletrônicos (CMOS x Nano). Por fim, foi realizada uma análise comparativa entre o consumo de energia das NoCs construídas com interconexões de cobre e das NoCs construídas com interconexões de BCNT, ambas com roteadores nanoeletrônicos.

5 - RESULTADOS E ANÁLISES

5.1 - INTRODUÇÃO

Neste capítulo, serão apresentadas as considerações adotadas e os resultados obtidos nas simulações e nos cálculos efetuados para realizar o estudo sobre o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos. As etapas descritas na metodologia foram seguidas, com o intuito de analisar o consumo total de energia dessa arquitetura de comunicação e comparar o consumo de energia entre NoCs com dispositivos nanoeletrônicos e NoCs com dispositivos CMOS. É importante ressaltar que em todos os casos estudados foram utilizados como base as seguintes considerações: topologia *mesh*, frequência de operação de 1 GHz, largura de canal com *flit* de 64 bits e interconexões entre os roteadores de 1 mm de comprimento, uma vez que esse valor é considerado razoável para representação de interconexões globais, pois está coerente com outros valores apresentados na literatura [3][14][66][78].

Ainda, para obtenção dos valores utilizados pelo modelo analítico de cálculo do consumo de energia, nesse trabalho foi adotado o modelo de rede sem contenção, onde o roteador necessita de três ciclos de *clock* para processar um *flit*, enquanto que a interconexão entre dois roteadores necessita de um ciclo para transmitir o *flit*. Este caso é semelhante ao roteador SCC de 48 núcleos da Intel sem contenção [79].

5.2 - OBTENÇÃO DO CONSUMO DE ENERGIA DAS INTERCONEXÕES

Nesta seção será realizada a análise das interconexões globais de cobre e BCNT de 1mm de comprimento, para vários nós de tecnologia (22 nm, 32 nm, 45 nm, 65 nm e 90 nm), utilizando dados disponibilizados por duas fontes diferentes, INTEL e ITRS.

Nos parâmetros disponibilizados pela INTEL, a sexta camada de interconexão foi escolhida como a camada responsável por interligar os roteadores dentro da NoC, pois foi considerado que a sexta, a sétima e a oitava camadas são identificadas como sendo aquelas responsáveis por realizar a roteamento global dentro da rede [79]. Para os parâmetros do ITRS, a largura utilizada na interconexão global foi de cinco vezes a largura mínima definida pelo nó de tecnologia para o fio de cobre, uma vez que a largura da interconexão global é tipicamente muito maior do que a largura mínima.

A Tabela 5.1 apresenta os parâmetros de largura, espessura e resistividade (ρ) para a interconexão global do material de cobre, em cada nó de tecnologia estudado, disponibilizados para os *chips* da INTEL [73][74][75][76][77] e pelo ITRS [33].

Tabela 5.1 – Parâmetros de interconexão retirados dos dados disponibilizados pela INTEL e dos relatórios do ITRS.

Nó de tecnologia	Parâmetros INTEL		Parâmetros ITRS		
	largura (nm)	espessura (nm)	largura (nm)	espessura (nm)	ρ ($\mu\Omega \cdot cm$)
22 nm	120	240	110	257	4,96
32 nm	169	303	160	374	4,08
45 nm	180	324	225	495	3,10
65 nm	240	430	325	715	2,73
90 nm	360	576	450	990	2,53

Como pode ser observado, os valores da resistividade do cobre utilizados nas camadas de interconexão da INTEL não estão disponíveis. Assim, a resistividade do cobre disponibilizada pelo ITRS também foi utilizada para obter a resistência do cobre para o modelo de interconexão construído a partir dos dados de interconexão da INTEL.

Após definidas as características da camada de interconexão global utilizada para interligar os roteadores da NoC, os parâmetros de resistência, indutância e capacitância foram obtidos para os modelos de cobre e BCNT (Tabelas B.1, B.2, B.3 e B.4 do Apêndice B). A constante dielétrica no valor de 2,6 foi utilizada para obter a capacitância do cobre, para todos os nós de tecnologia. Essa constante foi obtida a partir dos dados disponibilizados nos relatórios do ITRS [33].

Conforme visto anteriormente, a Equação (2.30) não reproduz corretamente o valor da capacitância eletroestática do BCNT. Neste sentido, com base nas observações e análises realizadas por Pasricha *et al.* [59], nesse trabalho foi considerado que o valor da capacitância total do BCNT por unidade de comprimento é igual ao do cobre, para a mesma seção transversal de interconexão. Ainda, assumiu-se na modelagem da camada de interconexão global que a largura do fio e o espaço entre as interconexões são iguais.

A partir da simulação da Figura 4-2, a frequência de corte em -3 dB e a potência para as interconexões globais foram obtidas, com a substituição do elemento de interconexão do esquemático, pelo circuito equivalente de cada nó de tecnologia estudado, tanto para o cobre, quanto para o BCNT. Assim, obteve-se a latência das interconexões, por meio da largura de banda do canal. A Figura 5-1 e a Figura 5-2 apresentam a latência das interconexões de cobre e de BCNT em função dos nós de tecnologia estudados, tanto para os parâmetros da INTEL, quanto para os parâmetros do ITRS, respectivamente.

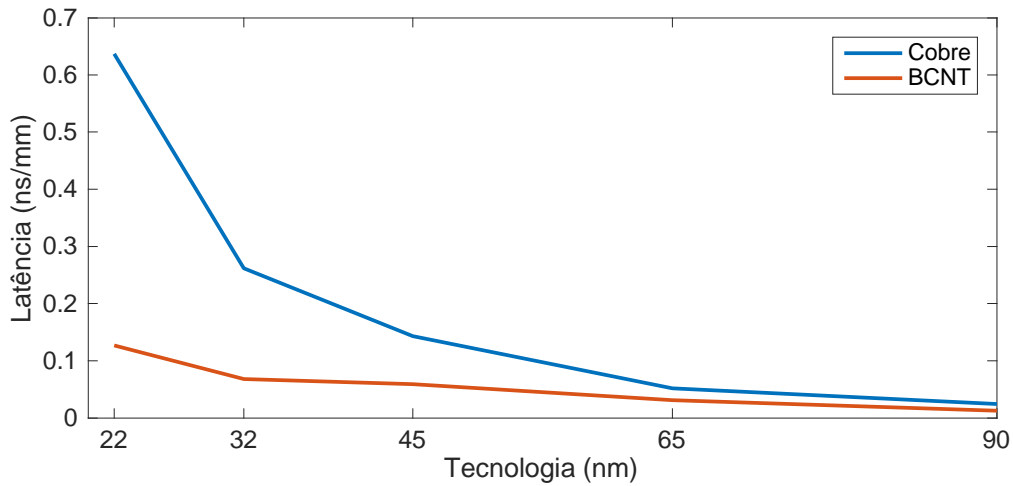


Figura 5-1 - Latência das interconexões de cobre e BCNT por nó de tecnologia para os dados disponibilizados pela INTEL.

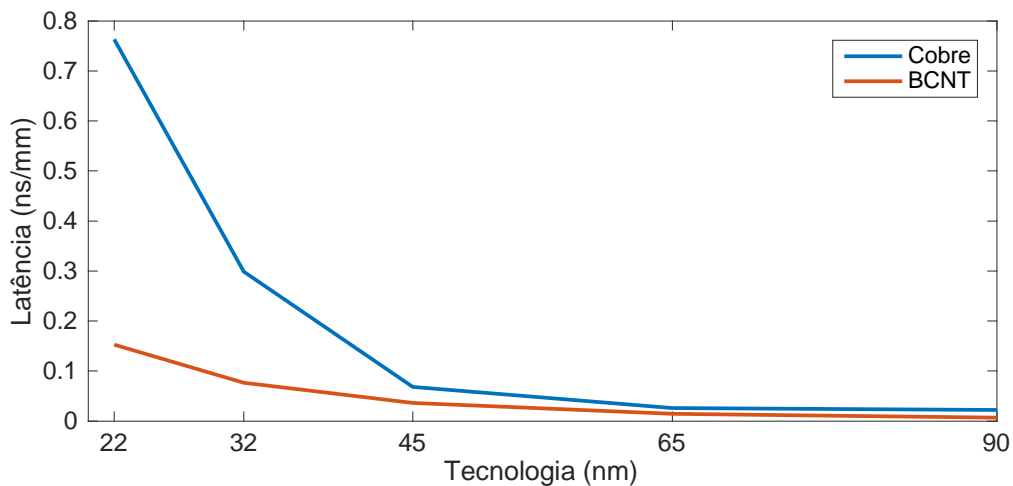


Figura 5-2 - Latência das interconexões de cobre e BCNT em função da tecnologia para os dados disponibilizados pelo ITRS.

Verifica-se em ambos os cenários que à medida que a tecnologia diminui, o atraso da interconexão aumenta, tanto para o cobre quanto para o BCNT. Isso ocorre devido à redução da largura e da espessura do fio que, por consequência, ocasiona o aumento da resistência. Além disso, verifica-se que a variação da latência é maior nas gerações tecnológicas de menores dimensões. Esse fato ocorre devido ao aumento da resistividade do cobre, ocasionado pela eletromigração e pelo aumento do espalhamento de contorno e superfície [44]. Ainda, na Figura 5-1 e na Figura 5-2 é mostrado que o BCNT possui melhor desempenho do que o cobre em todas as tecnologias. Assim, conclui-se que a utilização do material BCNT nas interconexões globais é bastante promissora para os futuros circuitos VLSI, conforme constatado em outros trabalhos [57][81][82][83].

Em seguida, obteve-se a energia por bit para os materiais de cobre e BCNT, por meio da simulação do circuito da Figura 4-2. A Figura 5-3 e a Figura 5-4 mostram o estudo comparativo da energia consumida por *bit* entre as interconexões de cobre e BCNT em função da tecnologia, para os parâmetros da INTEL e do ITRS.

Conforme pode ser visto na Figura 5-3 e na Figura 5-4, a energia consumida pela interconexão em todos os nós de tecnologia estudados é aproximadamente igual. Além disso, verifica-se que a energia consumida pela interconexão de BCNT é um pouco superior à do cobre. Vale ressaltar que para o modelo de interconexão simulado, não foi considerado o uso de repetidores. Assim, esse resultado é válido apenas para o cenário adotado sem o uso de repetidores. Ainda, pode ser observado que os valores da energia consumida em ambos os cenários são praticamente iguais. Isso ocorre, pois, a capacitância total obtida para os dois modelos de interconexão de cobre e a frequência de operação estão muito próximas. Assim, os resultados obtidos também são bastante semelhantes.

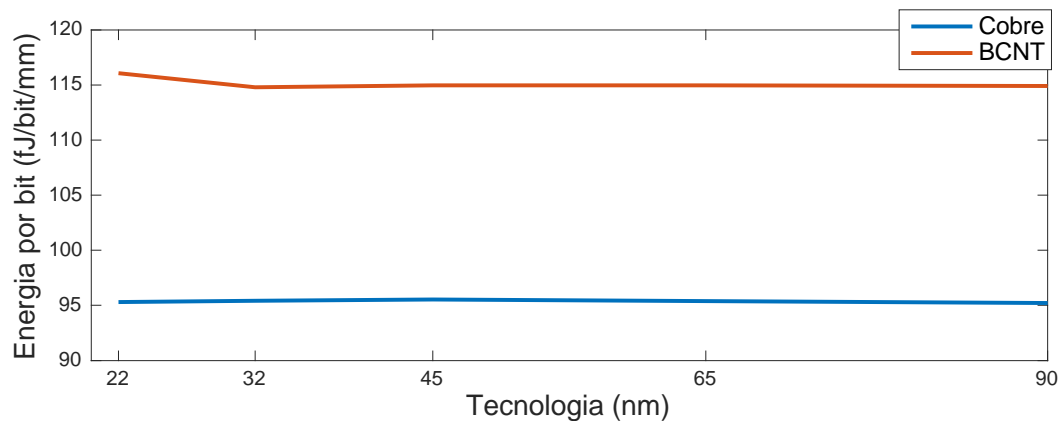


Figura 5-3 - Energia por bit em função do nó de tecnologia para os materiais de cobre e BCNT, parâmetros INTEL.

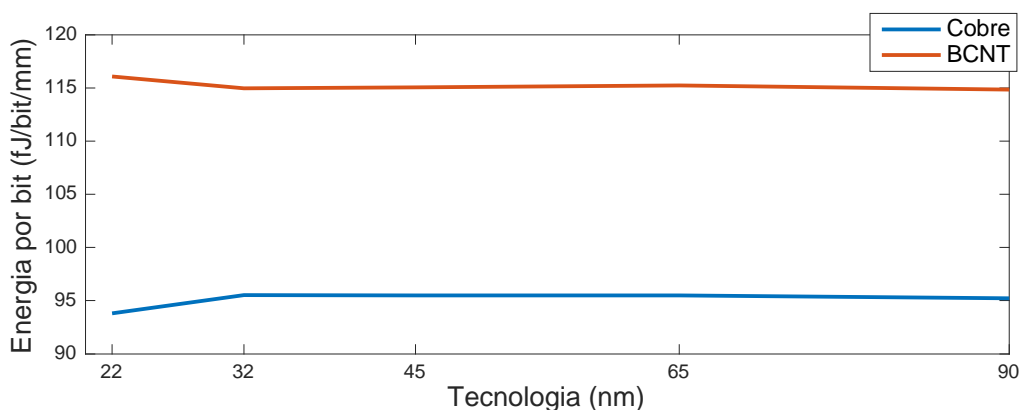


Figura 5-4 - Energia por bit em função do nó de tecnologia para os materiais de cobre e BCNT, parâmetros ITRS.

5.3 -OBTENÇÃO DO CONSUMO DE ENERGIA DO ROTEADOR

O valor total da potência dissipada pelo roteador nanoeletrônico, mostrado na Tabela 3.1, representa o consumo total desse dispositivo considerando todas as suas entradas e saídas [16]. No entanto, para utilizar o modelo de consumo de energia adotado, primeiramente foi necessário encontrar a energia consumida pela arquitetura do roteador nanoeletrônico, para processar apenas um *flit*. Dessa forma, por meio dos parâmetros da Tabela 3.1 e a partir da Equação (3.3) foi possível encontrar a energia consumida para processar um *bit* da entrada até a saída do roteador nanoeletrônico.

A arquitetura do roteador nanoeletrônico [16] foi realizada para tratamento de apenas um *bit* por entrada, utilizando um registrador que converte uma informação serial em paralela, para tratamento de uma palavra de 8 *bits*. Com o objetivo de calcular o consumo de energia do roteador nanoeletrônico para diferentes tamanhos de *flit*, e considerando que uma interconexão de uma NoC que chega em um roteador é composta por vários fios, faz-se necessário extrapolar o cálculo de potência fornecido por Câmara *et. al.* [16], dado que os componentes eletrônicos utilizados na construção desse roteador foram construídos para o tratamento de uma interconexão de apenas um fio, semelhante a uma interconexão composta por um *flit* de 1 bit.

Uma vez que a maior parte do consumo de energia do roteador é devida aos *buffers* e ao *crossbar* e dado que o consumo de energia desses componentes é proporcional ao tamanho do *flit* [69], nesta dissertação considera-se que o consumo de energia de cada roteador é diretamente proporcional ao tamanho do *flit*, ou seja, proporcional ao número de *bits* que o roteador processa por unidade de tempo. Conforme apresentado anteriormente, foi considerado que o roteador necessita de três ciclos de *clock* para processar um *flit*. Com base nessas considerações foi possível obter a energia consumida pelo roteador nanoeletrônico, para diferentes tamanhos de *flit*.

A fim de comparar a energia consumida para processar um *flit* pelo roteador nanoeletrônico e pelo roteador CMOS, a mesma metodologia apresentada anteriormente para calcular a energia consumida pelo roteador nanoeletrônico foi empregada para encontrar a energia consumida pelo roteador CMOS. A Figura 5-5 mostra a variação do consumo de energia do roteador nanoeletrônico e CMOS para diferentes tamanhos de *flit*, utilizando a extrapolação proposta e as considerações de retardo do modelo de rede adotado.

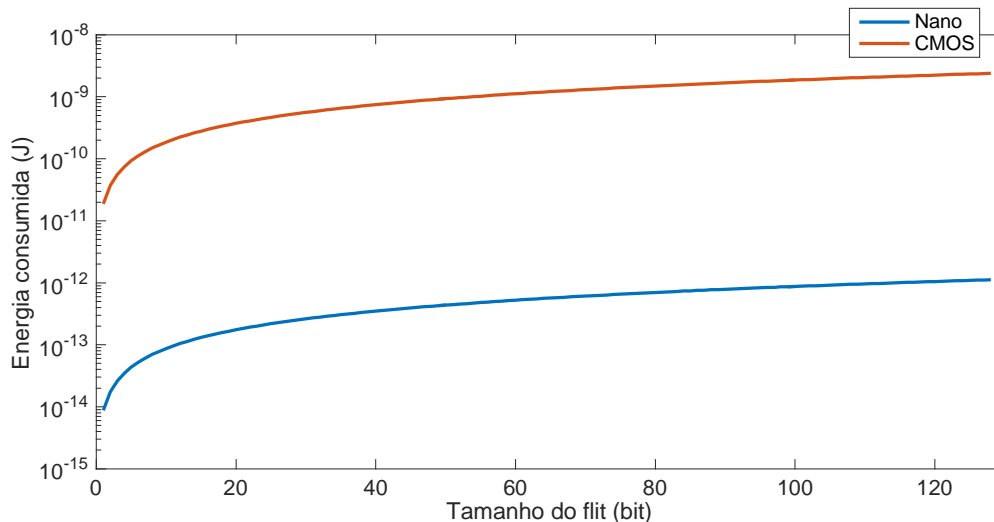


Figura 5-5 - Comparativo de energia consumida pelo roteador Nano e pelo roteador CMOS para diferentes tamanhos de flit.

Como pode ser visto, a energia consumida pelo roteador nanoeletrônico é bem inferior à do roteador CMOS e, portanto, aparenta ser uma proposta promissora para redução do consumo de energia dos futuros *chips* e dispositivos.

5.4 - OBTENÇÃO DO CONSUMO DE ENERGIA DE NoCs

Nesta seção, algumas NoCs são avaliadas no intuito de estudar sobre o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos. Para fins de comparação, as seguintes características da NoC foram alteradas: material empregado na interconexão e tecnologia empregada na construção do roteador. Assim, após seguir a metodologia proposta para obter o consumo de energia das interconexões de cobre e BCNT, bem como, o consumo de energia dos roteadores nanoeletrônicos e CMOS, o modelo analítico do consumo de energia escolhido [17] foi utilizado para calcular o total de energia consumida pela arquitetura de comunicação do tipo NoC, por meio da Equação (3.8).

Para utilizar essa equação, foi necessário definir o número de pacotes e o número de *flits* por pacote. Assim, para cada padrão de tráfego utilizado, 20.000 pacotes foram inseridos na rede, onde cada pacote possui 5 *flits* e cada *flit* possui 64 *bits* [17]. A energia total consumida pela NoC foi obtida pela soma da energia média de um *flit* sobre todas as distâncias de comunicação que ligam os nós dentro da topologia de rede escolhida, ponderadas pela CPD que depende do padrão de tráfego escolhido. Além disso, para o estudo sobre o consumo de energia de NoCs, foram utilizados os parâmetros disponibilizados pelo ITRS, dado que a maioria dos trabalhos científicos utiliza essa fonte e que a variação encontrada entre os cenários não foi muito significativa. A seguir será apresentado o consumo de energia e análises obtidos, para cada NoC estudada.

5.4.1 - Análise do Consumo de Energia de NoCs com Dispositivos Nanoeletrônicos

Para analisar o consumo de energia de NoCs com dispositivos nanoeletrônicos, primeiramente o modelo do consumo de energia adotado foi aplicado para uma NoC constituída por interconexões de cobre e por roteadores nanoeletrônicos. Adicionalmente, foi utilizada a topologia *mesh* e o algoritmo de roteamento XY para realizar os cálculos da CPD.

As constantes de energia da interconexão e do roteador para o *flit* foram obtidas por meio de simulações e estudos apresentados nas seções anteriores. Assim, para a tecnologia de 22 nm, o valor encontrado para o consumo de energia da interconexão global de cobre de 1 mm foi de aproximadamente 94 fJ/bit. Considerando que o *flit* adotado possui 64 bits e usando as considerações descritas na seção 3.1, o valor total do consumo de energia da interconexão encontrado foi de 6,016 pJ, enquanto que o valor encontrado para o consumo de energia utilizado pelo roteador nanoeletrônico para processar esse *flit* foi de 559,64 fJ.

A Figura 5-6 mostra o consumo de energia da NoC baseada em dispositivos nanoeletrônicos, para os diferentes padrões de tráfego apresentados anteriormente (uniforme, Rent, *Bit Complement*, *Bit Rotation*, *Nearest Neighbor*), em dois casos distintos, ideal e base. No caso ideal, apenas a energia dos roteadores é considerada, enquanto que no cenário base, tanto a energia dos roteadores, quanto a energia das interconexões são consideradas. Para o padrão de tráfego da regra de Rent, foi utilizando o expoente p no valor de 0,75, e para o padrão de tráfego *Nearest Neighbor* foi utilizado um raio de valor igual a 1 e fator de localidade de 50%.

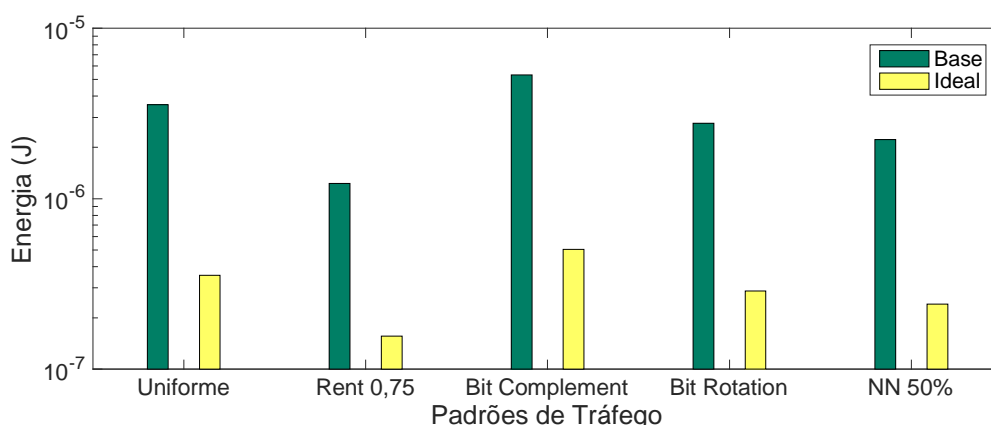


Figura 5-6 - Consumo total de energia de uma NoC 8X8 baseada em dispositivos Nanoeletrônicos.

Como pode ser observado na Figura 5-6, a energia dissipada pela comunicação da NoC baseada em dispositivos nanoeletrônicos é dominada pelo consumo de energia

das interconexões, em todos os padrões de tráfego estudados. Ainda, podemos observar no gráfico da Figura 5-6, que o menor consumo de energia é encontrado no padrão de tráfego da regra de Rent. Este resultado poderia ser pressuposto a partir das CPD da Figura 3-4 (e), visto que esse é o tráfego que possui maior localidade de comunicação. Além disso, como o tráfego da regra de Rent é baseado em dados experimentais, espera-se que esse tráfego forneça um melhor modelo de localidade de comunicação de aplicativos reais do que os demais padrões de tráfego sintéticos.

5.4.2 - Análise do Consumo de Energia de NoCs com dispositivos CMOS

Para fins de comparação, as condições aplicadas na seção anterior são utilizadas para analisar o consumo de energia de NoCs com dispositivos CMOS. Assim, o modelo é aplicado em uma NoC constituída por interconexões de cobre de 1 mm e roteadores implementados em tecnologia CMOS. Dado que o mesmo tipo de material foi utilizado para essa configuração de rede, o mesmo valor encontrado para o consumo de energia da interconexão é utilizado. A energia dissipada pelo roteador CMOS para processar o flit é de 1,196 nJ.

A Figura 5-7 mostra o consumo de energia da NoC baseada em dispositivos CMOS, para os mesmos padrões de tráfego e condições apresentadas anteriormente. Assim, do mesmo modo que foi feito para a NoC baseada em dispositivos nanoeletrônicos, dois casos distintos foram estudados, ideal e base.

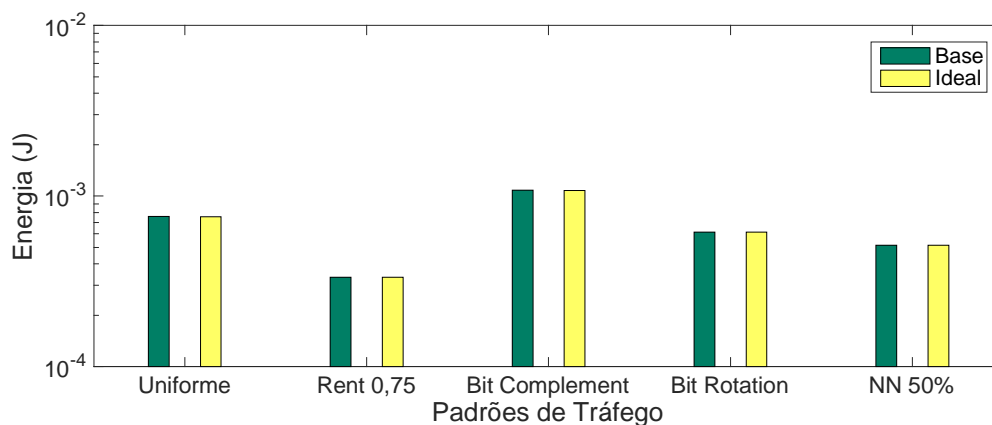


Figura 5-7 - Consumo total de energia de uma NoC 8X8 baseada em dispositivos CMOS.

Pode-se observar na Figura 5-7 que a maior parte da dissipação de energia de uma NoC baseada em dispositivos CMOS é ocasionada pelo consumo de energia dos roteadores CMOS, ao contrário do que ocorre na NoC baseada em dispositivos nanoeletrônicos. Portanto, como pode ser visto, em termos de energia, os roteadores são o principal gargalo da NoC baseada em dispositivos CMOS, enquanto que na NoC baseada em dispositivos nanoeletrônicos, as interconexões são o principal gargalo.

5.4.3 - Comparação do Consumo de Energia de NoCs CMOS X NoCs Nano

A Figura 5-8 mostra a comparação do consumo de energia entre NoCs com dispositivos CMOS e NoCs com dispositivos nanoeletrônicos, para vários tamanhos da topologia da NoC, do tipo *mesh* quadrada. Nessa análise, considerou-se o uso de interconexões globais do nó de tecnologia de 22 nm, *flit* de 64 bits e padrão de tráfego uniforme para calcular a CPD.

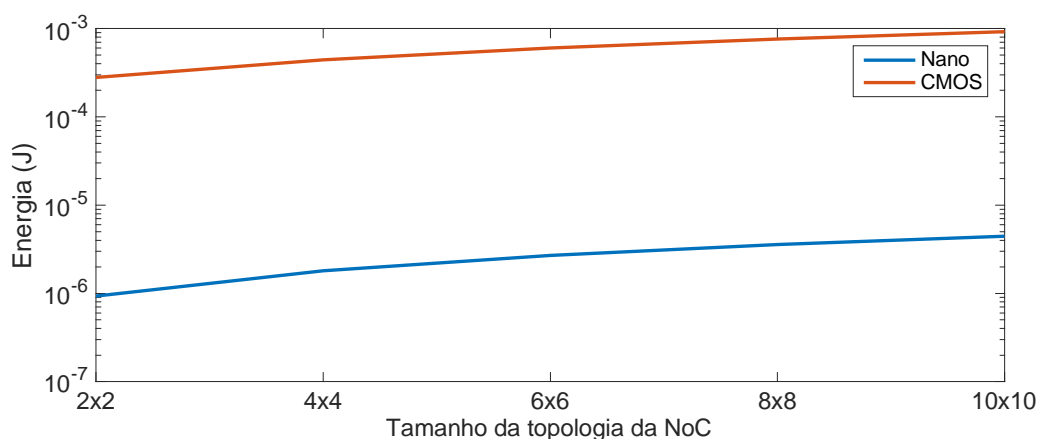


Figura 5-8 - Comparativo do consumo total de energia entre uma NoC baseada em dispositivos nanoeletrônicos e uma NoC baseada em dispositivos CMOS.

Conforme esperado, quanto maior o tamanho da topologia da NoC, maior será a quantidade de roteadores e interconexões dentro da rede e, conseqüentemente, maior o consumo de energia. Ainda, a rede baseada em dispositivos CMOS possui um consumo de energia bem superior ao das redes baseadas em dispositivos nanoeletrônicos. Assim, conforme apresentado nos resultados obtidos, pode-se concluir que a tecnologia utilizada na construção dos roteadores é o principal contribuinte na energia consumida pela NoC.

5.4.4 - Comparação do Consumo de Energia entre NoCs com Interconexões de Cobre e BCNT

A Figura 5-9 apresenta a comparação do consumo de energia de NoCs construídas com interconexões de cobre e NoCs construídas com interconexões de BCNT, em função do tamanho da topologia da NoC, do tipo *mesh* quadrada. Nesse estudo, foram utilizados roteadores implementados com tecnologia nanoeletrônica, interconexões globais do nó de tecnologia de 22 nm e padrão de tráfego uniforme para calcular a CPD.

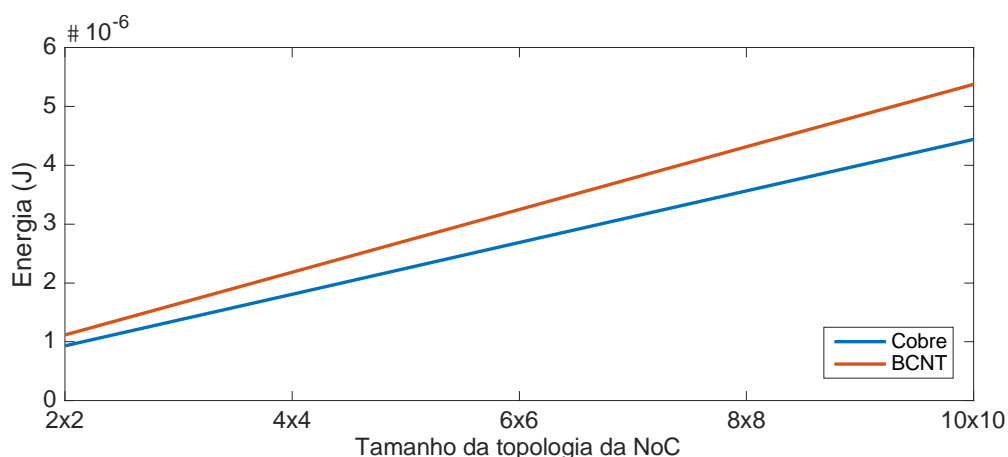


Figura 5-9 - Comparativo do consumo total de energia entre uma NoC construída com interconexões de cobre e uma NoC construída com interconexões de BCNT, em função do tamanho da topologia da NoC.

Conforme esperado, com o aumento do tamanho da topologia da NoC, a energia consumida por essa arquitetura de comunicação também aumenta. Além disso, percebe-se que a variação da energia consumida pela NoC com interconexões de BCNT é maior do que para a NoC com interconexões de cobre. Isso ocorre, pois, a contribuição da energia das interconexões é a parcela predominante no cálculo da energia total consumida pela NoC com dispositivos nanoeletrônicos e, com isso, para o modelo de interconexão sem repetidor, a interconexão de BCNT consome mais energia do que a de cobre. Assim, com o aumento da topologia da NoC, o crescimento do consumo de energia da rede que utiliza interconexões BCNT tende a ser maior do que o da rede com interconexões de cobre.

A Figura 5-10 mostra outro cenário para a comparação do consumo de energia de NoCs construídas com interconexões de cobre e NoCs construídas com interconexões de BCNT. Neste cenário, o consumo de energia é avaliado em função do nó de tecnologia. Para essa NoC estudada, a topologia utilizada foi uma malha 8x8 e o tráfego utilizado no cálculo da CPD foi o padrão uniforme.

Como pode ser observado na Figura 5-10, é possível verificar que a variação de energia da NoC é semelhante a observada na Figura 5-4, em concordância com os resultados obtidos anteriormente, onde as interconexões são a parcela predominante da energia total consumida pela NoC baseada em dispositivos nanoeletrônicos.

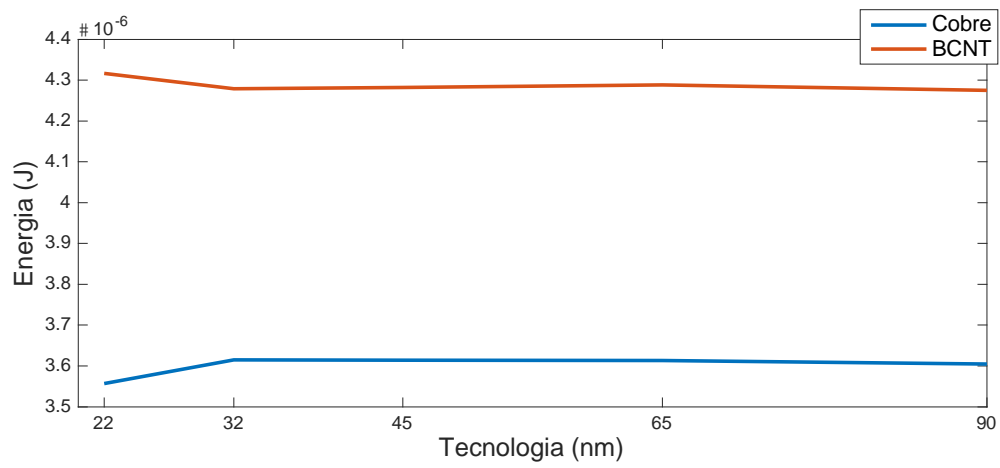


Figura 5-10 - Comparativo do consumo total de energia de uma NoC construída com interconexões de cobre e outra construída com interconexões de BCNT, em função do nó de tecnologia.

6 - CONCLUSÕES E PERSPECTIVAS FUTURAS

Neste trabalho, foi realizado o estudo sobre o consumo de energia em NoCs baseadas em dispositivos nanoeletrônicos, utilizando o modelo analítico proposto por Bezerra [17], que calcula o consumo de energia de NoCs, por meio da sua CPD. O modelo adotado assume que o consumo de energia é proporcional à distância percorrida pelos pacotes dentro da rede. Dado que uma NoC é constituída basicamente por interconexões e roteadores, definiu-se realizar a avaliação do desempenho das NoCs, por meio da substituição da tecnologia empregada nesses dois componentes, onde para as interconexões foram utilizados e avaliados os materiais de cobre e BCNT, e para os roteadores foram empregadas as tecnologias CMOS e nanoeletrônica.

As interconexões globais de cobre e BCNT foram simuladas e analisadas, considerando as dimensões da camada responsável por interligar os roteadores dentro da rede e utilizando dados disponibilizados por duas fontes diferentes, INTEL e ITRS. Verificou-se que com a redução das dimensões dos dispositivos e da tecnologia, o atraso da interconexão aumenta, tanto para o cobre, quanto para o BCNT. Além disso, foi verificado que as interconexões globais de BCNT possuem melhor desempenho do que as de cobre. Com relação ao consumo de energia, foi observado que a interconexão de BCNT consome mais energia do que a interconexão de cobre, para modelos sem repetidor. Conforme apresentado, visando atingir menores valores de latência nas interconexões, o emprego de outros materiais, como o BCNT, pode ser visto como uma possível solução para a continuidade da redução do número de ciclos que interligam dois núcleos dentro de uma NoC.

Em seguida, utilizando a arquitetura do roteador nanoeletrônico baseado em SETs, a energia consumida para o roteador da NoC processar um *flit* foi calculada, nas duas tecnologias estudadas. Assim, verificou-se que a energia consumida pelo roteador nanoeletrônico é bem inferior à do roteador CMOS.

Após encontrados os valores de energia consumida tanto para as interconexões de cobre e BCNT, quanto para os roteadores CMOS e nanoeletrônico, o consumo de energia da NoC com roteadores nanoeletrônicos foi encontrado, em cinco padrões de tráfego diferentes. Com isso, foi verificado que a energia dissipada pela comunicação baseada em dispositivos nanoeletrônicos é dominada pela energia dissipada pelas interconexões. Para fins de comparação, o consumo de energia de uma NoC com dispositivos CMOS foi encontrado. Para essa configuração, verificou-se que a energia dissipada pela NoC é dominada pela energia consumida pelos roteadores, ao contrário da NoC com dispositivos nanoeletrônicos. Além disso, obteve-se o comparativo do consumo de energia da NoC, em função do tamanho da rede. Conforme apresentado, verificou-se que a tecnologia utilizada na construção dos roteadores é o principal contribuinte da energia consumida pela NoC.

Por fim, foi realizada a comparação do consumo de energia entre NoCs construídas com interconexões de cobre e NoCs construídas com interconexões de BCNT, ambas com roteadores nanoeletrônicos. Nessa análise, foram considerados dois casos, onde o primeiro avaliou o consumo de energia em função do tamanho da topologia da rede e o segundo avaliou o consumo de energia em função do nó de tecnologia. Em ambos os casos, foi verificado que o consumo de energia da NoC segue a variação de energia imposta pelas interconexões.

A comunicação em *chip* consome uma parcela significativa de potência e área do *chip*. Assim, visando reduzir a energia dissipada em uma NoC, a partir do estudo realizado nesse trabalho, verificou-se quantitativamente que o uso de roteadores nanoeletrônicos, aparenta ser uma proposta promissora para a redução do consumo de energia total da NoC e, portanto, conclui-se que a nanoeletrônica é uma tecnologia que apresenta ser uma solução para reduzir o consumo de energia dos futuros *chips* e dispositivos.

Os resultados obtidos neste trabalho advêm dos parâmetros escolhidos para as interconexões, bem como dos critérios adotados para a rede, tais como, falta de contenção, roteamento XY e topologia *mesh*. Ainda, vale ressaltar que o uso de repetidores não foi considerado nos modelos de interconexão utilizados nas NoCs estudadas. Neste sentido, como perspectivas futuras, é importante investigar um modelo de interconexão com parâmetros otimizados que utilize repetidores, dado que estes também contribuem com latência e dissipação de energia, e com isso, aprimorar a arquitetura de comunicação baseada em dispositivos nanoeletrônicos proposta neste trabalho, possibilitando a sua comparação com outras NoCs reais. Além disso, sugere-se estudar o consumo de energia de NoCs baseadas em dispositivos nanoeletrônicos, em outras topologias, para averiguar qual seria aquela onde a tecnologia nanoeletrônica poderia se sobressair. Por fim, sugere-se aplicar a NoC baseada em dispositivos nanoeletrônicos em um sistema completo, para analisar o desempenho total de uma aplicação, com a utilização da tecnologia nanoeletrônica, e averiguar a parcela de contribuição de energia dessa arquitetura de comunicação, no consumo total de energia do *chip*.

REFERÊNCIAS BIBLIOGRÁFICAS

- [1] MOORE, G. E. Cramming more components onto integrated circuits. **Proceedings of the IEEE**, v. 86, n. 1, p. 82-85, 1998.
- [2] DENNARD, R. H. Design of ion-implanted MOSFET's with very small physical dimensions. **IEEE Journal of Solid-State Circuits**, v. 9, n. 5, p. 256-268, 1974.
- [3] KRISHNA, T. (2014). **Enabling dedicated single-cycle connections over a shared network-on-chip**. Tese de Doutorado, Massachusetts Institute of Technology.
- [4] VENKATESH, G. Conservation cores: reducing the energy of mature computations. In: **ACM SIGARCH Computer Architecture News**. ACM, p. 205-218, 2010.
- [5] DAVIS, J. A. Interconnect limits on gigascale integration (GSI) in the 21st century. **Proceedings of the IEEE**, v. 89, n. 3, p. 305-324, 2001
- [6] SRIVASTAVA A.; XU, Y.; SHARMA, A. K. Carbon nanotubes for next generation very large scale integration interconnects. **Journal of Nanophotonics**, v. 4, n. 041690, p. 1-26, 2010.
- [7] SRIVASTAVA N.; LI, H.; KREUPL, F.; BANERJEE, K. On the Applicability of Single-Walled Carbon Nanotubes as VLSI Interconnects. **IEEE Transactions on Nanotechnology**, v. 8, n. 4, p. 542-558, 2009.
- [8] THIRUVENKATESAN, C.; RAJA, J. Studies on the Application of Carbon Nanotube as Interconnects for Nanometric VLSI Circuits. **ICETET'09**, p. 162-167, 2009.
- [9] DAS, D.; RAHAMAN, H. Timing Analysis in Carbon Nanotube Interconnects with Process, Temperature, and Voltage Variations. **2010 International Symposium on Electronic System Design**, p. 27-32, 2010.
- [10] BHAT, S. Energy models for network-on-chip components. **Master of Science, Department of Mathematics and Computer Science, Technische Universiteit Eindhoven, Eindhoven**, 2005.

- [11] PATTERSON, D. A. **Computer architecture: a quantitative approach**. Elsevier, 2011.
- [12] FOCHI, V. M. (2015). **Técnicas de tolerância a falhas aplicadas a redes intra-chip**. Tese de Doutorado, Pontifícia Universidade Católica do Rio Grande do Sul.
- [13] WANG, H.; PEH, L. S.; MALIK, S. Power-driven design of router microarchitectures in on-chip networks. In: **Proceedings of the 36th annual IEEE/ACM International Symposium on Microarchitecture**. IEEE Computer Society, 2003. p. 105.
- [14] VANGAL, S. R. (2007). **Performance and Energy Efficient Network-on-Chip Architectures**. Tese de Doutorado, Linköping University. Linköping, Suécia, 93 p.
- [15] PALMA, J. C. S. (2007). **Reduzindo o Consumo de Potência em Networks-on-Chip através de Esquemas de Codificação de Dados**. Tese de Doutorado. Programa de Pós-Graduação em Computação. UFRGS, Porto Alegre.
- [16] CÂMARA, B. O. (2017). **Roteador Nanoeletrônico para Redes-em-Chip baseado em Transistores Monoelétron**. Dissertação de Mestrado, Publicação 656/2017 DM/PGEA, Departamento de Engenharia Elétrica, Universidade de Brasília, Brasília, DF, 75 p.
- [17] BEZERRA, G. B. (2012) **Energy consumption in networks on chip: efficiency and scaling**. Tese de Doutorado, The University of New Mexico.
- [18] SCHALLER, R. Moore's law: past, present and future. **IEEE Spectrum**, v.34, n.6, 1997, pp. 53-59.
- [19] ZEFERINO, C. A. (2003). **Redes-em-Chip: Arquiteturas e Modelos para Avaliação de Área e Desempenho**. Tese de Doutorado, Porto Alegre: UFRG, 242 p.
- [20] CHEN, W.T.; SHEU, J.P. Performance Analysis of Multiple Bus Interconnection Networks with Hierarchical Requesting Model. **IEEE Transactions on Computers**, v. 40, n. 7, pp. 834- 842, 1991.

- [21] IBM. The CoreConnect Bus Architecture. White Paper, 1999. 8 p.
- [22] BENINI, L.; DE MICHELI, G. Networks on Chips: a New SoC Paradigm. **IEEE Computer**, v.35, n.1, 2002, pp. 70-78.
- [23] DALLY, W.; TOWLES, B. P. Route Packets, not Wires: On-Chip Interconnection Networks. **In: Design Automation Conference (DAC'01)**, 2001, pp. 684-689.
- [24] DALLY, W. J.; TOWLES, B. P. **Principles and practices of interconnection networks**. Elsevier, 2004.
- [25] DUATO, J.; YALAMANCHILI, S.; NI, L. M. **Interconnection networks: an engineering approach**. Morgan Kaufmann, 2003.
- [26] PÉREZ, A. R. (2012). **Floorplan-Aware High Performance NoC Design**. Tese de Doutorado, Universitat Politècnica de València. València, Espanha.
- [27] OULD-KHAOUA, M.; MIN, G. Circuit Switching: An Analysis for k-Ary n-Cubes with Virtual Channels. **IEE Proceedings – Computers and Digital Techniques**, v. 148, n. 6, pp. 215-219, 2001.
- [28] DALLY, W.J.; SEITZ, C.L. Deadlock-Free Message Routing in Multiprocessor Interconnection Networks. **IEEE Transactions on Computers**, v. 36, n. 5, pp. 547-553, 1987.
- [29] PEH, L.S.; DALLY, W.J. Flit-Reservation Flow Control. **In: Proceedings of 6th International Symposium on High-Performance Computer Architecture (HPCA'00)**, Toulouse, pp. 73-84, 2000.
- [30] DALLY, W.J. Virtual Channel Flow Control. **IEEE Transactions on Parallel and Distributed Systems**, v. 3, n. 2, pp. 194-205, 1992.
- [31] TEDESCO, L. P. **Uma Proposta para Geração de Tráfego e Avaliação de Desempenho para NoCs**. 2005. Tese de Doutorado. Pontifícia Universidade Católica do Rio Grande do Sul.

- [32] NICOPOULOS, C. A. **Network-on-Chip architectures: A holistic design exploration**. The Pennsylvania State University, 2007.
- [33] International Technology Roadmap for Semiconductors (ITRS). 2016. Disponível em: <http://www.itrs2.net/itrs-reports/>.
- [34] REEHAL, G. (2012) **Designing Low Power and High Performance Network-on-Chip Communication Architectures for Nanometer SoCs**. Tese de Doutorado. The Ohio State University.
- [35] CEYHAN, A. (2014). **Interconnects for future technology generations-conventional CMOS with copper/low-k and beyond**. Tese de Doutorado. Georgia Institute of Technology.
- [36] MEINDL, J. D. Beyond Moore's law: The interconnect era. **Computing in Science & Engineering**, v. 5, n. 1, p. 20-24, 2003.
- [37] JAN, M. R.; ANANTHA, C.; BORIVOJE, N. Digital Integrated Circuits—A Design Perspective. 2003.
- [38] HO, R.; MAI, K. W.; HOROWITZ, M. A. The future of wires. **Proceedings of the IEEE**, v. 89, n. 4, p. 490-504, 2001.
- [39] EDELSTEIN, D. Full copper wiring in a sub-0.25/ μm CMOS ULSI technology. In: **Electron Devices Meeting, 1997. IEDM'97. Technical Digest, International**. IEEE, 1997. p. 773-776.
- [40] BOHR, M. The new era of scaling in an SoC world. In: **Solid-State Circuits Conference-Digest of Technical Papers, 2009. ISSCC 2009. IEEE International**. IEEE, 2009. p. 23-28.
- [41] NAEEMI, A.; SARVARI, R.; MEINDL, J. D. Performance comparison between carbon nanotube and copper interconnects for gigascale integration (GSI). **IEEE Electron Device Letters**, v. 26, n. 2, p. 84-86, 2005.

- [42] SINGH, D. P.; RAI, M. K. G. (2013). **Study The Waveform Analysis of Interconnect Performance For Future Vlsi Design**. Dissertação de mestrado. Thapar University, Punjab, Índia.
- [43] MAFFUCCI, A. Electrical Conductivity of Carbon Nanotubes: Modeling and Characterization. In: **Carbon Nanotubes for Interconnects**. Springer International Publishing, 2017. p. 101-128.
- [44] KOO, K. H. (2011). **The Comparison Study of Future On-chip Interconnects for High Performance VLSI Applications**. Tese de Doutorado. Stanford University.
- [45] BAI, P. A 65nm logic technology featuring 35nm gate lengths, enhanced channel strain, 8 Cu interconnect layers, low-k ILD and 0.57/spl mu/m/sup 2/SRAM cell. In: **Electron Devices Meeting, 2004. IEDM Technical Digest. IEEE International**. IEEE, 2004. p. 657-660.
- [46] RAI, M. K.; SARKAR, S. Influence of tube diameter on C nanotube interconnect delay and power output. **Physica Satus Solidi A**, v. 298, p. 735-739.
- [47] SRIVASTAVA, N.; BANERJEE, K. A comparative scaling analysis of metallic and carbon nanotube interconnections for nanometer scale VLSI technologies. In: **Proc. 21st Intl. VLSI Multilevel Interconnect Conf**. 2004. p. 393-398.
- [48] JANTSCH, A.; TENHUNEN, H. (Eds.) **Networks on chip**. Dordrecht: Kluwer Academic Publishers, 2003.
- [49] BAKOGLU, H. B. **Circuits, Interconnections, and Packaging for VLSI**. 1990.
- [50] NOGUEIRA, C. P. S. M. (2012). **Análise Comparativa entre Interconexões de Nanotubo de Carbono e Interconexões de Cobre para Circuitos GSI/TSI**. Dissertação de Mestrado em Engenharia de Sistemas Eletrônicos e de Automação, Publicação PPGEA.DM-488/2012, Departamento de Engenharia Elétrica, Universidade de Brasília, Brasília, DF, 72p.

- [51] KURUVILLA, N.; RAINA, J. P. Statistical latency analysis of carbon nanotube interconnects due to contact resistance variations. In: **Microelectronics, 2008. ICM 2008. International Conference on**. IEEE, 2008. p. 296-299.
- [52] NAEEMI, A.; MEINDL, J. D. Design and performance modeling for single-walled carbon nanotubes as local, semiglobal, and global interconnects in gigascale integrated systems. **IEEE Transactions on Electron Devices**, v. 54, n. 1, p. 26-37, 2007.
- [53] RAHAMAN, M. S.; CHOWDHURY, M. H. Information theoretic capacity analysis of single-walled carbon nanotube bundle VLSI interconnects. In: **Integrated Circuits, ISIC'09. Proceedings of the 2009 12th International Symposium on**. IEEE, 2009. p. 530-533.
- [54] LI, H. Circuit modeling and performance analysis of multi-walled carbon nanotube interconnects. **IEEE Transactions on electron devices**, v. 55, n. 6, p. 1328-1337, 2008.
- [55] NIEUWOUDT, A.; MASSOUD, Y. On the impact of process variations for carbon nanotube bundles for VLSI interconnect. **IEEE Transactions on Electron Devices**, v. 54, n. 3, p. 446-455, 2007.
- [56] BURKE, Peter J. Luttinger liquid theory as a model of the gigahertz electrical properties of carbon nanotubes. **IEEE Transactions on Nanotechnology**, v. 99, n. 3, p. 129-144, 2002.
- [57] SRIVASTAVA, N.; BANERJEE, K. Performance analysis of carbon nanotube interconnects for VLSI applications. In: **Proceedings of the 2005 IEEE/ACM International conference on Computer-aided design**. IEEE Computer Society, 2005. p. 383-390.
- [58] XU, Y.; SRIVASTAVA, A.; SHARMA, A. K. Emerging carbon nanotube electronic circuits, modeling, and performance. **VLSI Design**, v. 2010, p. 7, 2010.

- [59] PASRICHA, S.; DUTT, N.; KURDAHI, F. J.; Exploring carbon nanotube bundle global interconnects for chip multiprocessor applications. In: **VLSI Design, 2009 22nd International Conference on**. IEEE, 2009. p. 499-504.
- [60] UCHINO, T.; CONG, J. An interconnect energy model considering coupling effects. **IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems**, v. 21, n. 7, p. 763-776, 2002.
- [61] LINEAR. Linear Technology - Design Simulation and Device Models. 2017. Disponível em: <http://www.linear.com/designtools/software>.
- [62] PANDE, P. P.; GRECU, C.; JONES M.; IVANOV A.; SALEH R. Effect of traffic localization on energy dissipation in NoC-based interconnect. In: **Circuits and Systems, 2005. ISCAS 2005. IEEE International Symposium on**. IEEE, 2005. p. 1774-1777.
- [63] CHRISTIE, P.; STROOBANDT, D. The interpretation and application of Rent's rule. **IEEE Transactions on Very Large Scale Integration (VLSI) Systems**, v. 8, n. 6, p. 639-648, 2000.
- [64] DAVIS, J. A.; DE, V. K.; MEINDL, J. D. A stochastic wire-length distribution for gigascale integration (GSI) - Part I: Derivation and validation. **IEEE Transactions on Electron Devices**, v. 45, n. 3, p. 580-589, 1998.
- [65] STROOBANDT, D. **A priori wire length estimates for digital design**. Kluwer Academic Publishers, Boston, 2001
- [66] HOSKOTE, Y. A 5-GHz mesh interconnect for a teraflops processor. **IEEE Micro**, v. 27, n. 5, p. 51-61, 2007.
- [67] TAYLOR, M. B.; KIM J.; MILLER J.; WENTZLAFF D.; GHODRAT F.; GREENWALD B.; HOFFMAN H.; JOHNSON P.; LEE W. The raw microprocessor: A computational fabric for software circuits and general-purpose programs. **IEEE micro**, v. 22, n. 2, p. 25-35, 2002.

- [68] HEIRMAN, W. DAMBRE, J.; STROOBANDT, D.; CAMPENHOUT, J. V. Rent's rule and parallel programs: characterizing network traffic behavior. In: **Proceedings of the 2008 international workshop on System level interconnect prediction**. ACM, 2008. p. 87-94.
- [69] KAHNG, A.; LI, B.; PEH, L.S.; SAMADI, K. ORION 2.0: a fast and accurate NoC power and area model for early-stage design space exploration. In: **Proceedings of the conference on Design, Automation and Test in Europe**. European Design and Automation Association, 2009. p. 423-428.
- [70] HU, J.; MARCULESCU, R. Energy-aware mapping for tile-based NoC architectures under performance constraints. In: **Proceedings of the 2003 Asia and South Pacific Design Automation Conference**. ACM, 2003. p. 233-239.
- [71] PALMA, J. C. S.; MARCON, C. A. M.; MORAES, F.G.; CALAZANS, N.L.V.; REIS, R.A.L.; SUSIN, A.A. Mapping embedded systems onto NoCs: the traffic effect on dynamic energy estimation. In: **Proceedings of the 18th annual symposium on Integrated circuits and system design**. ACM, 2005. p. 196-201.
- [72] PANDE, P. P.; GRECU, C.; JONE, M.; IVANOV, A.; SALEH, R. Effect of traffic localization on energy dissipation in NoC-based interconnect. In: **Circuits and Systems, 2005. ISCAS 2005. IEEE International Symposium on**. IEEE, 2005. p. 1774-1777.
- [73] INGERLY, D. *et al.* Low-k interconnect stack with metal-insulator-metal capacitors for 22nm high volume manufacturing. In: **Interconnect Technology Conference (IITC), 2012 IEEE International**. IEEE, 2012. p. 1-3
- [74] PACKAN, P. *et al.*; High Performance 32nm Logic Technology Featuring 2 nd Generation High-k+ Metal Gate Transistors. In: **Electron Devices Meeting (IEDM), 2009 IEEE International**. IEEE, 2009. p. 1-4.
- [75] MOON, P. *et al.*; Process and Electrical Results for the On-die Interconnect Stack for Intel's 45nm Process Generation. **Intel Technology Journal**, v. 12, n. 2, 2008.

- [76] BAI, P. *et al.*; A 65nm logic technology featuring 35nm gate lengths, enhanced channel strain, 8 Cu interconnect layers, low-k ILD and 0.57/spl mu/m/sup 2/SRAM cell. In: **Electron Devices Meeting, 2004. IEDM Technical Digest. IEEE International**. IEEE, 2004. p. 657-660
- [77] JAN, C.-H. *et al.*; 90 nm generation, 300 mm wafer low k ILD/Cu interconnect technology. In: **Interconnect Technology Conference, 2003. Proceedings of the IEEE 2003 International**. IEEE, 2003. p. 15-17.
- [78] CHEN, G. *et al.* A 340 mV-to-0.9 V 20.2 Tb/s source-synchronous hybrid packet/circuit-switched 16× 16 network-on-chip in 22 nm tri-gate CMOS. **IEEE Journal of Solid-State Circuits**, v. 50, n. 1, p. 59-67, 2015.
- [79] HOWARD, J. *et al.*; A 48-core IA-32 message-passing processor with DVFS in 45nm CMOS. In: **Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2010 IEEE International**. IEEE, 2010. p. 108-109.
- [80] JAN, C.H. *et al.*; A 22nm SoC platform technology featuring 3-D tri-gate and high-k/metal gate, optimized for ultra low power, high performance and high density SoC applications. In: **Electron Devices Meeting (IEDM), 2012 IEEE International**. IEEE, 2012. p. 3.1. 1-3.1. 4.
- [81] CHEN, F. *et al.* Scaling and evaluation of carbon nanotube interconnects for VLSI applications. In: **Proceedings of the 2nd international conference on Nano-Networks**. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2007. p. 24.
- [82] LI, H.; YIN, W. Y.; MAO, J. F.; Modeling of carbon nanotube interconnects and comparative analysis with Cu interconnects. In: **Microwave Conference, 2006. APMC 2006. Asia-Pacific**. IEEE, 2006. p. 1361-1364
- [83] JOSHI, A.; SONI, G.; A Comparative Analysis of Copper and Carbon Nanotubes-Based Global Interconnects in 32 nm Technology. In: **Proceedings of Fifth International Conference on Soft Computing for Problem Solving**. Springer, Singapore, 2016. p. 425-437.

- [84] LIENTSCHNIG, G.; WEYMANN, I.; HADLEY, P.; Simulating hybrid circuits of single-electron transistors and field-effect transistors. **Japanese journal of applied physics**, v. 42, n. 10R, p. 6467, 200
- [85] KARIMIAN, M. *et al.*; A new SPICE macro-model for simulation of single electron circuits. In: **Microelectronics (ICM), 2009 International Conference on**. IEEE, 2009. p. 228-231
- [86] WU, Y. L.; LIN, S. T.; An improved single-electron-transistor model for spice application. **Nanotechnology**, v. 3, p. 321, 2003.
- [87] GUIMARÃES, J. G. (2005). “**Arquiteturas de Redes Neurais Nanoeletrônicas para Processadores em escala Giga ou Tera**”. Tese de Doutorado, Universidade de Brasília. Brasília, Brasil.

APÊNDICES

A - FUNCIONAMENTO DO SET

O transistor mono-elétron (SET) consiste de duas junções túnel conectadas em série, formando uma ilha entre elas, como mostrado na Figura A-1 [84][85][86].

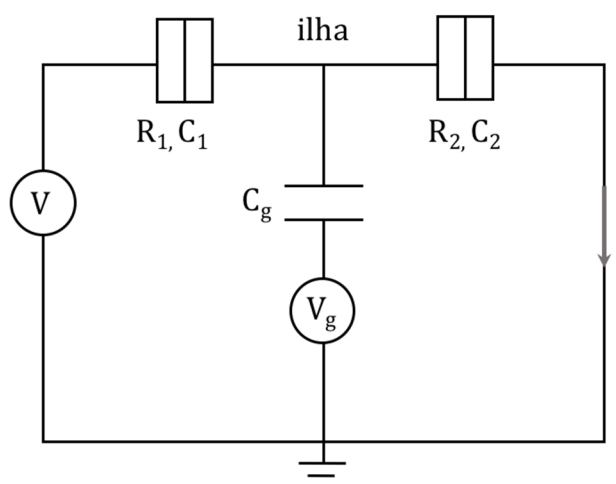


Figura A-1 - Transistor mono-elétron.

As junções túnel são constituídas por dois eletrodos metálicos separados por um isolante muito fino (barreira), o qual permite a passagem de elétrons por tunelamento [84][85][86]. Os parâmetros R_1 , C_1 e R_2 , C_2 da Figura A-1 correspondem às resistências e capacitâncias da primeira e segunda junção túnel, respectivamente, que dependem da área e da espessura da barreira isolante. A tensão de porta V_g controla a energia eletroestática da ilha por meio da capacitância C_g . Assim, quando há um carregamento por tunelamento de uma junção e descarregamento da outra junção, o fluxo controlado de cargas gera a corrente I [87]. Dessa forma, pode-se resumir que a operação do SET é baseada no controle de fluxo de elétrons individuais por tunelamento, o que conduz a um consumo de energia muito baixo.

B - TABELAS COMPLEMENTARES

Tabela B.1- Parâmetros obtidos a partir do modelo da interconexão global de cobre para a fonte INTEL.

Tecnologia	R_{Cu} (Ω)	L_{Cu} (H)	C_{Cu} (F)
22 nm	1,72E+03	3,57E-09	1,42E-13
32 nm	7,99E+02	3,45E-09	1,41E-13
45 nm	5,32E+02	3,42E-09	1,41E-13
65 nm	2,65E+02	3,31E-09	1,41E-13
90 nm	1,22E+02	3,16E-09	1,41E-13

Tabela B.2 – Parâmetros obtidos a partir do modelo da interconexão global de cobre para a fonte ITRS.

Tecnologia	R_{Cu} (Ω)	L_{Cu} (H)	C_{Cu} (F)
22 nm	2,05E+03	3,60E-09	1,42E-13
32 nm	8,85E+02	3,47E-09	1,41E-13
45 nm	3,40E+02	3,33E-09	1,41E-13
65 nm	1,44E+02	3,18E-09	1,41E-13
90 nm	7,81E+01	3,07E-09	1,41E-13

Tabela B.3 – Parâmetros obtidos a partir do modelo da interconexão global de BCNT para a fonte INTEL.

Tecnologia	R_Q (Ω)	R_S (Ω)	L_{CNT} (H)
22 nm	1,71E-01	1,14E+02	7,08E-11
32 nm	9,62E-02	6,41E+01	3,98E-11
45 nm	8,42E-02	5,62E+01	3,48E-11
65 nm	4,75E-02	3,17E+01	1,96E-11
90 nm	2,36E-02	1,57E+01	9,75E-12

Tabela B.4 – Parâmetros obtidos a partir do modelo da interconexão global de BCNT para a fonte ITRS.

Tecnologia	R_Q (Ω)	R_S (Ω)	L_{CNT} (H)
22 nm	2,04E-01	1,36E+02	8,44E-11
32 nm	1,07E-01	7,11E+01	4,41E-11
45 nm	5,38E-02	3,59E+01	2,22E-11
65 nm	2,57E-02	1,71E+01	1,06E-11
90 nm	1,51E-02	1,00E+01	6,23E-12