

筑波大学計算科学センターシンポジウム
「計算機アーキテクトが考える
次世代スパコン」

2006年4月5日

村上和彰

九州大学

murakami@cc.kyushu-u.ac.jp

次世代スパコン

～達成目標と制約条件の整理～

- 達成目標

- 性能目標(2011年)

- LINPACK (HPL): 10PFlop/s
- 実アプリケーション: 1PFlop/s

- 成果目標

- (私見) 科学技術計算能力の国際競争力の向上ならびに維持による我が国の科学技術力、産業力の発揚

- 制約条件

次世代スパコン

～達成目標と制約条件の整理～

- 達成目標
- 制約条件
 - Scalability: NLSとしての存在のみならず、NISへの下方展開が可能なこと
 - Productivity: ソフトウェアの生産性が高いこと
 - Compatibility: 過去のソフトウェア資産からの、ならびに、将来のソフトウェア開発への連続性、継続性を保証すること
 - Dependability: 長期間連続運転可能なこと
 - Serviceability: センター運用可能なこと
 - Ecology: 低消費電力
 - Economy: 将来にわたってビジネス的に持続的に開発可能であること

10PFlop/sのマシンをどう作るか？

クロック周波数？
演算器数？
プロセッサ数？

メモリバンド幅？
メモリレイテンシ？
メモリサイズ？

LINPACK: 10PFlop/s
↓ 効率60%を仮定
ピーク性能: 16PFlop/s

.....

計算ノード

プロセッサ

メモリ

計算ノード

プロセッサ

メモリ

計算ノード

プロセッサ

メモリ

システムインターコネク

p2p通信バンド幅？
p2p通信レイテンシ？
バイセクションバンド幅？

10PFlop/sのマシンをどう作るか？

達成目標性能: LINPACK 10PFlop/s

↓ 仮定: 実行効率60%

ピーク性能: 16PFlop/s

仮定: クロック周波数1GHz

演算器数: 16M個

仮定: 4-way MUL&ADD SIMD

プロセッサコア数: 2M個

仮定: 4-way CMP

プロセッサチップ数: 512K個

プロセッサチップ性能: 32GFlop/s

仮定: 4B/s@Flop/s

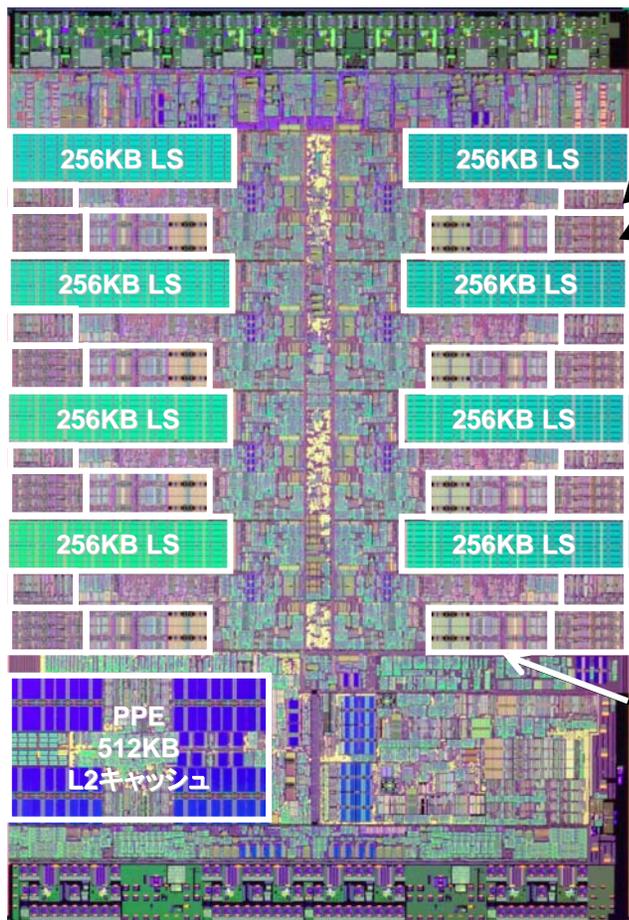
要求オフチップメモリバンド幅: 128GB/s

現在のプロセッサチップは押し並べて...

i -way superscalar processor with j -way SIMD × k -way CMP

	i	j	k	クロック 周波数	プロセッサ チップ性能	オフチップメモリ バンド幅
Intel Pentium	3	4	2	4GHz	64GFlop/s	
IBM BG/L	2	4	2	800MHz	12.8GFlop/s	
Cell	2	4	8	4GHz	256GFlop/s	25.6GB/s
ClearSpeed	1	96	1	200MHz	38.4GFlop/s	

現在のプロセッサチップは押し並べて...



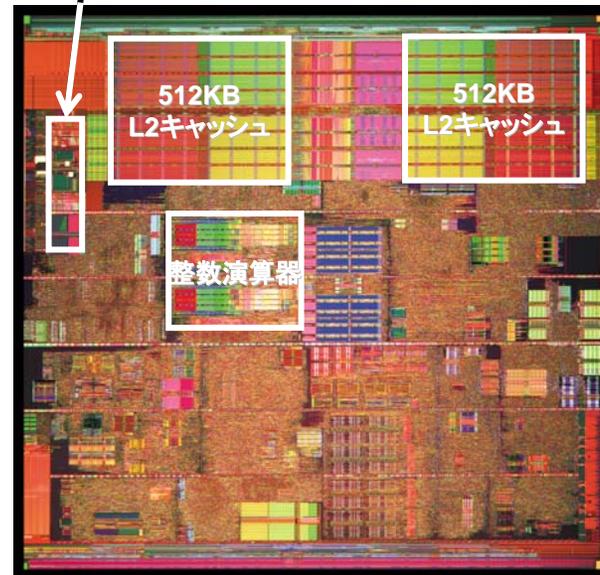
Cell (221mm² @90nm)

倍精度浮動小数点演算器

単精度浮動小数点演算器

浮動小数点演算器

整数演算器

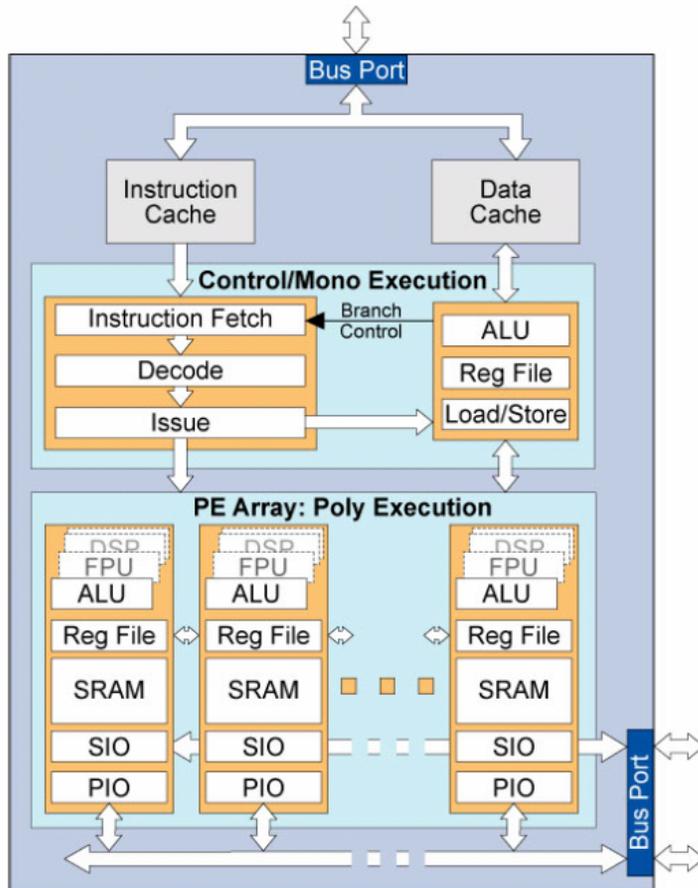


Pentium4 (112mm² @90nm)

CellとPentium4の比較

	CELL	Pentium4 (Prescott)
製造技術	90nm SOI	90nm
トランジスタ数	2億3400万	1億2500万
面積	221mm ²	112mm ²
動作周波数	4.0GHz(最大4.6GHz)	3.8GHz
ピーク性能	256GFlop/s(SPE × 8個)	15.2GFlop/s
消費電力	48W(6W × 8SPE) + α	103W
プロセッサコア	9個	1個
メインプロセッサ	PPE(Powerベース+VMX) × 1個	IA32(x86)ベース(+SSE3)
演算プロセッサ	SPE(SIMDプロセッサ) × 8個	—
内部メモリ		
メインプロセッサ	<ul style="list-style-type: none"> •L1命令キャッシュ:32KB •L1データキャッシュ:32KB •L2キャッシュ:512KB 	<ul style="list-style-type: none"> •L1命令キャッシュ:12Kマイクロ命令 •L1データキャッシュ:16KB •L2キャッシュ:1MB
演算プロセッサ	•各SPEのLS:256KB	—
チップ内インターコネクト (オンチップバス)	EIB:192Gバイト/秒((128ビット+64ビット) ×2GHz×4リング)	—
外付けDRAMインタフェース	XIO:25.6Gバイト/秒(32ビット×3.2GHz×2 チャンネル)	FSB:6.4Gバイト/秒(64ビット ×800MHz)
入出カインタフェース	FlexIO:76.8Gバイト/秒(6.4Gビット/秒×8ビット ×12トランシーバ/レシーバ)	

Multi-Threaded Array Processing Architecture



- Multi-Threaded Array Processor
 - Fully programmable in C
 - Hardware multi-threading
 - Extensible instruction set
- Scalable internal parallelism
 - Array of processing elements (PEs)
 - Compute, bandwidth scale together
 - From 10s to 1,000s of PEs
 - Built-in PE redundancy
- High performance, low power
 - ~10 GFLOPS/watt
- Multiple high-speed I/O channels

10PFlop/sのマシンをどう作るか？

達成目標性能: LINPACK 10PFlop/s

↓ 仮定: 実行効率60%

ピーク性能: 16PFlop/s

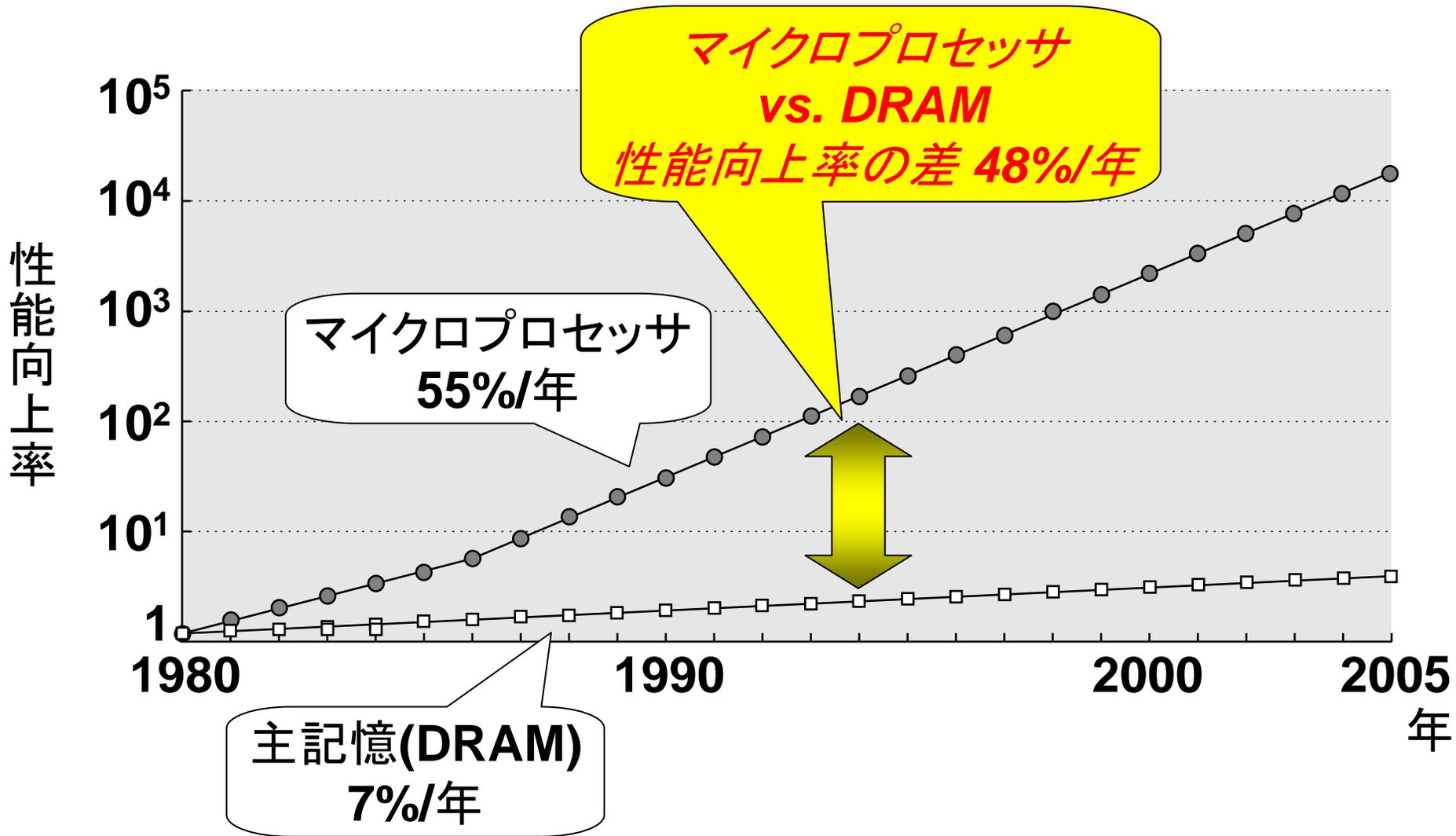
仮定: FB-DIMM

チャンネル当りメモリバンド幅: 8GB/s

仮定: 4B/s@Flop/s

要求オフチップメモリバンド幅: 128GB/s

「Memory Wall」問題



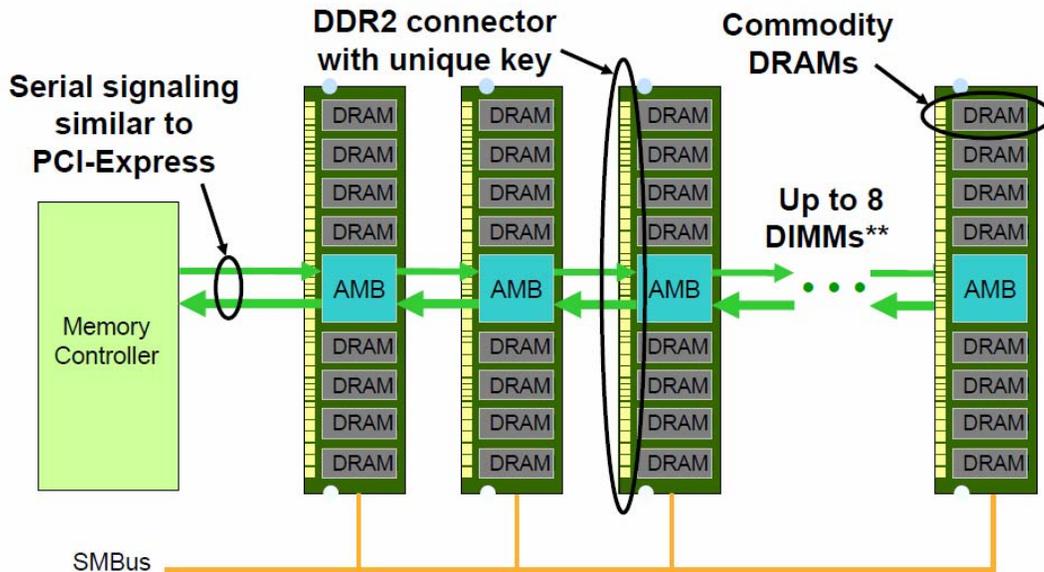
DDR/DDR2/DDR3 SDRAM

項目	DDR3 SDRAM	DDR2 SDRAM	DDR SDRAM
クロック周波数	533/667MHz	200/266/333/400MHz	100/133/166/200MHz
データ転送速度	1066/1333Mbps	400/533/667/800Mbps	200/266/333/400Mbps
I/O幅	x4/x8/x16	x4/x8/x16	x4/x8/x16/x32
プリフェッチ	8ビット	4ビット	2ビット
クロック入力	ディファレンシャルクロック	ディファレンシャルクロック	ディファレンシャルクロック
バースト長	4, 8, 4 (バーストチョップ)	4, 8	2, 4, 8
データストロープ形式	ディファレンシャルデータストロープ	ディファレンシャルデータストロープ	シングルデータストロープ
電源電圧	1.5V	1.8V	2.5V
I/Oインタフェース	SSTL_15	SSTL_18	SSTL_2
/CASレーテンシ(CL)	5, 6, 7, 8, 9, 10クロック	3, 4, 5クロック	2, 2.5, 3クロック
On die termination (ODT)	対応	対応	非対応
単体パッケージ	FBGA	FBGA	TSOP(II)/FBGA/LQFP
鉛フリー	対応	対応	対応

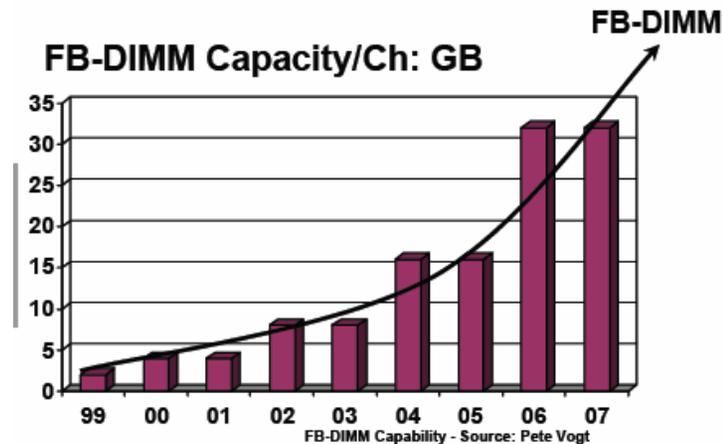
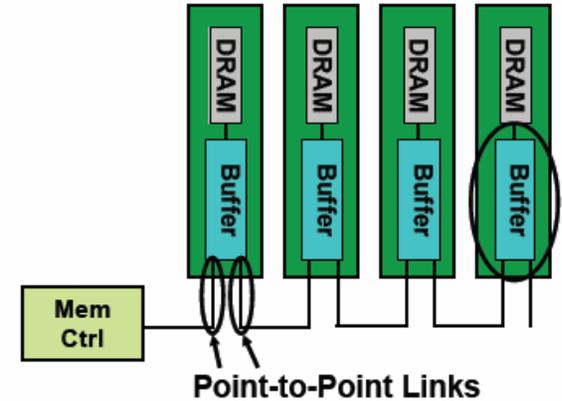
<http://www.elpida.com/pdfs/J0876E10.pdf>

FB-DIMM

- アドレスやデータ等のバッファリング
- P2P通信
- 総メモリ容量のスケールリングを可能に！
(最大8DIMM)



** FBD spec supports up to 8 DIMMs per channel, however initial Intel chipsets will only address 4 DIMMs.



FB-DIMM製品例(エルピーダ)

- 2005年8月サンプル出荷

4GB FB-DIMM 主な仕様

Density	4GB	
Organization	256M words x 72-bits x 2 ranks	
Module Speed Grade	PC2-4200F	PC2-5300F
DRAM Speed Grade	DDR2-533 (533Mbps)	DDR2-667 (667Mbps)
FB-DIMM Channel Peak Throughput	6.4GB/s	8.0GB/s
Mounted Devices	1Gb DDR2 SDRAM x 36 pieces	
Burst Length	4 or 8	
Supply Voltage	1.8V±0.1V	
Package	240-pin Fully Buffered DIMM 133.35mm x 30.35mm Thin:6.7mm (max.)	



10PFlop/sのマシンをどう作るか？

達成目標性能: LINPACK 10PFlop/s

↓ 仮定: 実行効率60%

ピーク性能: 16PFlop/s

仮定: FB-DIMM

チャンネル当りメモリバンド幅: 8GB/s

仮定: 4B/s@Flop/s

要求オフチップメモリバンド幅: 128GB/s

「Memory Wall」問題を抱えて、
どうプロセッサを作るか？

10PFlop/sのマシンをどう作るか？

達成目標性能: LINPACK 10PFlop/s

↓ 仮定: 実行効率60%

ピーク性能: 16PFlop/s

仮定: クロック周波数1GHz

演算器数: 16M個

仮定: 4-way MUL&ADD SIMD

プロセッサコア数: 2M個

仮定: 4-way CMP

プロセッサチップ数: 512K個

プロセッサチップ性能: 32GFlop/s

仮定: 4B/s@Flop/s

要求オフチップメモリバンド幅: 128GB/s

仮定: クロック周波数1GHz

演算器数: 16M個

仮定: 1024-way RDP

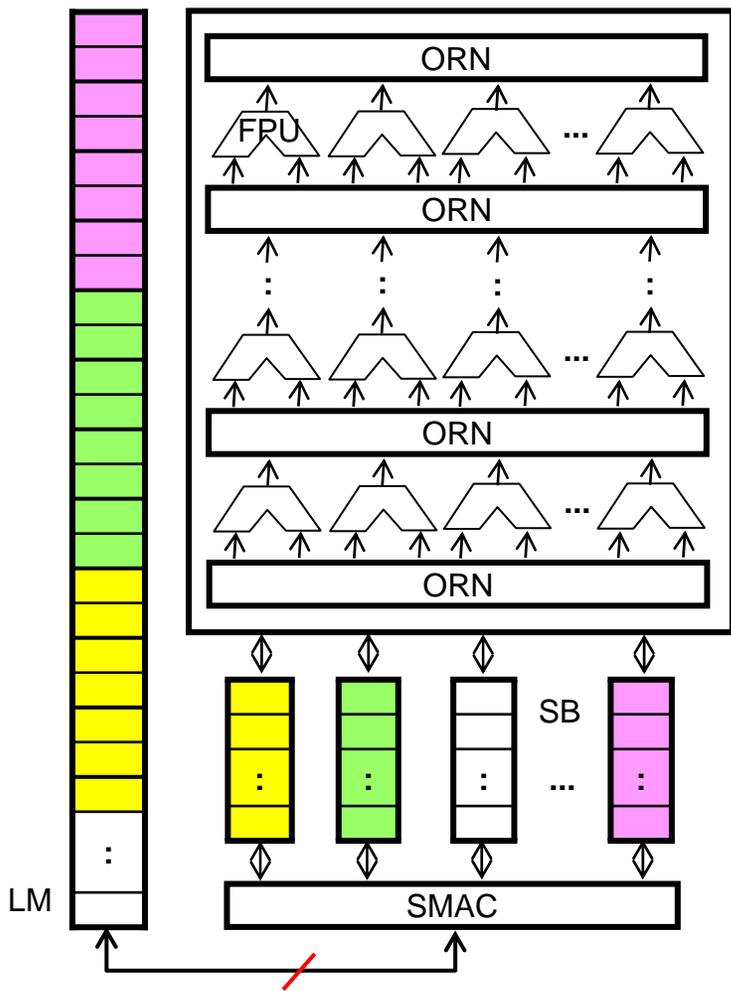
プロセッサチップ数: 16K個

プロセッサチップ性能: 1TFlop/s

仮定: 0.0625B/s@Flop/s

要求オフチップメモリバンド幅: 64GB/s

提案：再構成可能大規模データパス (RDP: Reconfigurable Datapath)



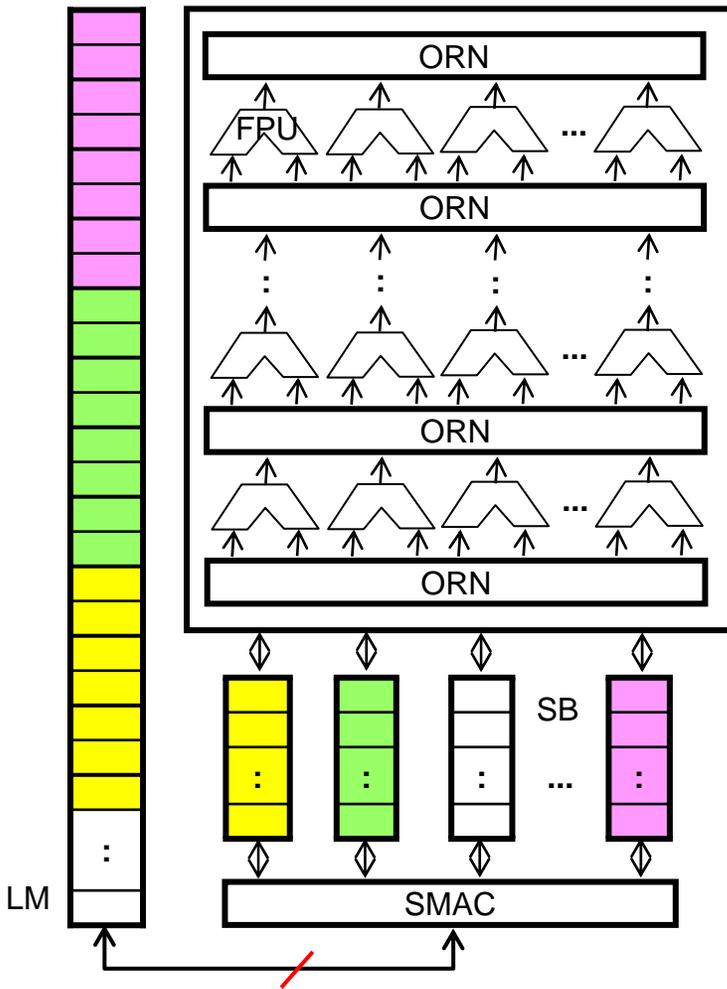
• RDPとは？

– 多数の演算器 (FPU: Floating-Point Unit) とそれらを相互接続する網 (ORN: Operand Routing Network) を搭載し、

- FPUで行う演算内容
- ORN上のFPU間接続関係

を動的に再構成可能としたデータパス

提案：再構成可能大規模データパス (RDP: Reconfigurable Datapath)



趣旨は？

— トランジスタ資源を(従来のプロセッサのように)、

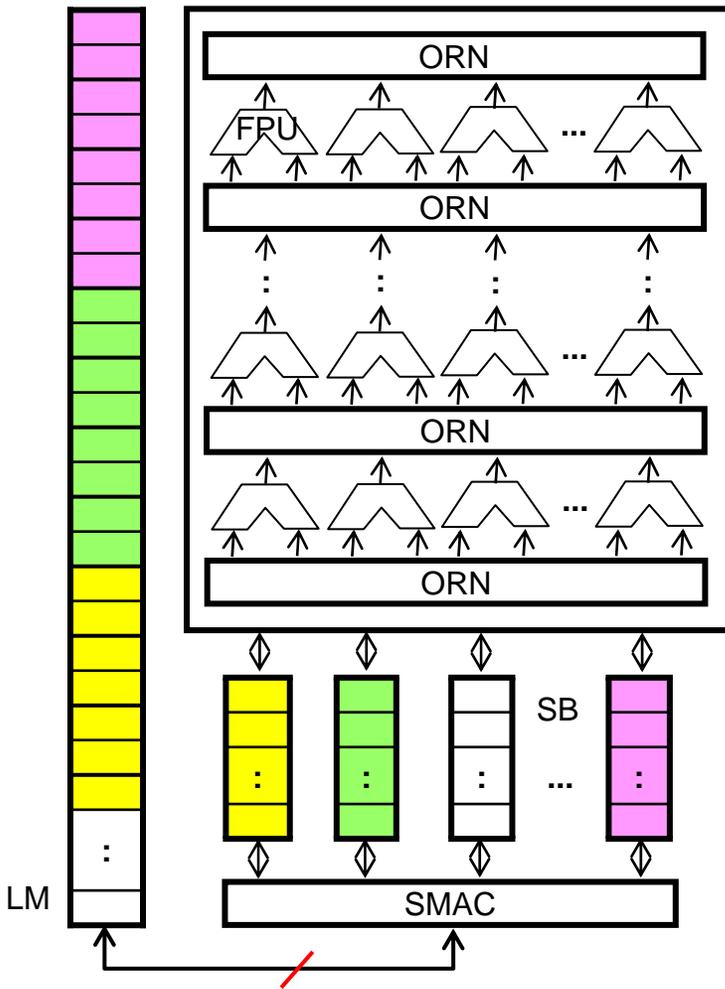
- データ並列性の活用(SIMD、ベクトルプロセッサ)
- 命令レベル並列性の活用(OOOスーパースカラ)
- プロセスレベル並列性の活用(CMP)

に投資するのではなく、

- データ依存性の維持(ORN上でのデータ転送)
- イタレーション間並列性の活用(複数のイタレーションをパイプライン処理)

に投資することで実効性能を確保！

提案：再構成可能大規模データパス (RDP: Reconfigurable Datapath)



• 用途は？

- 主プロセッサに対するコプロセッサ
- コア計算部のループボディのデータフロー全体を直接マッピング
 - ORNでデータ依存関係を維持
 - 原理的には毎クロックサイクル、新しいイタレーションを実行開始可能→複数イタレーションのパイプライン処理

RDPの応用例

～分子軌道法における二電子積分計算 ($\mu\nu \parallel \lambda\sigma$) の場合～

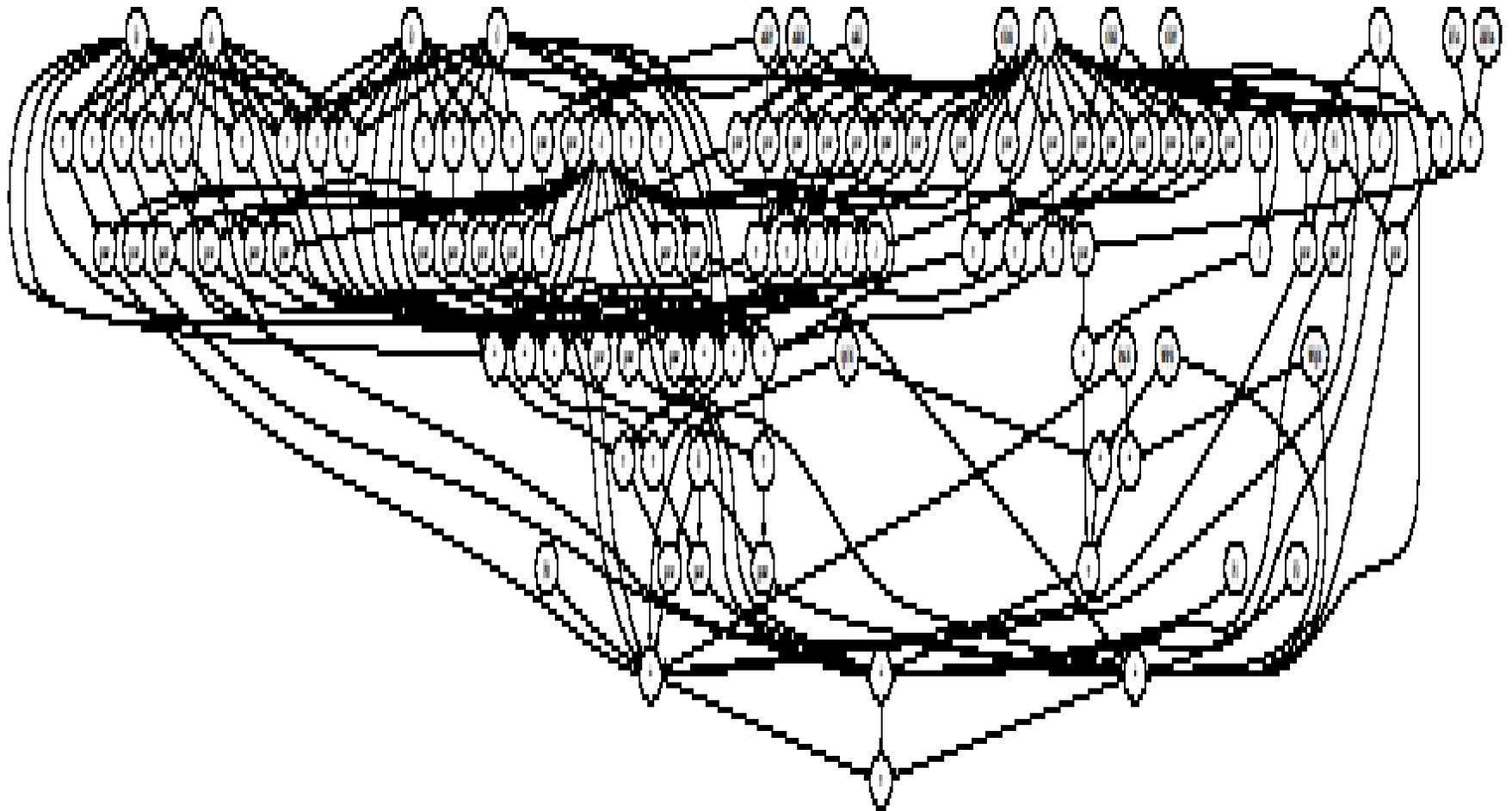
$$\text{tei}(4,4,4,4)=(((3+2*p*(4*PAx*PBx+PBx**2+PAx**2*(1+2*p*PBx**2)))*(3+2*q*(4*QCx*QDx+QDx**2+QCx**2*(1+2*q*QDx**2)))*f(0,t))/(p**2*q**2)+(4*(3+2*p*(4*PAx*PBx+PBx**2+PAx**2*(1+2*p*PBx**2)))*PQx*(QCx+QDx)*(3+2*q*QCx*QDx)*f(1,t))/(p*q*(p+q))*(4*(PAx+PBx)*(3+2*p*PAx*PBx)*PQx*(3+2*q*(4*QCx*QDx+QDx**2+QCx**2*(1+2*q*QDx**2)))*f(1,t))/(p*q*(p+q))*(8*(PAx+PBx)*(3+2*p*PAx*PBx)*(QCx+QDx)*(3+2*q*QCx*QDx)*((p+q)*f(1,t)+2*p*PQx**2*q*f(2,t)))/(p*q*(p+q)**2)+(2*(3+2*p*(4*PAx*PBx+PBx**2+PAx**2*(1+2*p*PBx**2)))*(3+q*(QCx**2+4*QCx*QDx+QDx**2))*((p+q)*f(1,t)+2*p*PQx**2*q*f(2,t)))/(p*q**2*(p+q)**2)+(2*(3+p*(PAx**2+4*PAx*PBx+PBx**2)))*(3+2*q*(4*QCx*QDx+QDx**2+QCx**2*(1+2*q*QDx**2))*((p+q)*f(1,t)+2*p*PQx**2*q*f(2,t)))/(p*q**2*(p+q)**2)+(4*(3+2*p*(4*PAx*PBx+PBx**2+PAx**2*(1+2*p*PBx**2)))*PQx*(QCx+QDx)*(3*(p+q)*f(2,t)+2*p*PQx**2*q*f(3,t)))/(q*(p+q)**3)+(8*(3+p*(PAx**2+4*PAx*PBx+PBx**2))*PQx*(QCx+QDx)*(3+2*q*QCx*QDx)*(3*(p+q)*f(2,t)+2*p*PQx**2*q*f(3,t)))/(p*(p+q)**3)*(8*(PAx+PBx)*(3+2*p*PAx*PBx)*PQx*(3+q*(QCx**2+4*QCx*QDx+QDx**2))*((p+q)*f(2,t)+2*p*PQx**2*q*f(3,t)))/(q*(p+q)**3)*(4*(PAx+PBx)*PQx*(3+2*q*(4*QCx*QDx+QDx**2+QCx**2*(1+2*q*QDx**2)))*(3*(p+q)*f(2,t)+2*p*PQx**2*q*f(3,t)))/(p*(p+q)**3)+((3+2*p*(4*PAx*PBx+PBx**2+PAx**2*(1+2*p*PBx**2)))*(3*(p+q)**2*f(2,t)+4*p*PQx**2*q*(3*(p+q)*f(3,t)+p*PQx**2*q*f(4,t)))/(q**2*(p+q)**4)*(8*(PAx+PBx)*(3+2*p*PAx*PBx)*(QCx+QDx)*(3*(p+q)**2*f(2,t)+4*p*PQx**2*q*(3*(p+q)*f(3,t)+p*PQx**2*q*f(4,t)))/(q*(p+q)**4)*(8*(PAx+PBx)*(QCx+QDx)*(3+2*q*QCx*QDx)*(3*(p+q)**2*f(2,t)+4*p*PQx**2*q*(3*(p+q)*f(3,t)+p*PQx**2*q*f(4,t)))/(p*(p+q)**4)+(4*(3+p*(PAx**2+4*PAx*PBx+PBx**2)))*(3+q*(QCx**2+4*QCx*QDx+QDx**2))*((3*(p+q)**2*f(2,t)+4*p*PQx**2*q*(3*(p+q)*f(3,t)+p*PQx**2*q*f(4,t)))/(p*q*(p+q)**4)+((3+2*q*(4*QCx*QDx+QDx**2+QCx**2*(1+2*q*QDx**2)))*(3*(p+q)**2*f(2,t)+4*p*PQx**2*q*(3*(p+q)*f(3,t)+p*PQx**2*q*f(4,t)))/(p**2*(p+q)**4)*(4*p*(PAx+PBx)*(3+2*p*PAx*PBx)*PQx*(15*(p+q)**2*f(3,t)+4*p*PQx**2*q*(5*(p+q)*f(4,t)+p*PQx**2*q*f(5,t)))/(q*(p+q)**5)+(8*(3+p*(PAx**2+4*PAx*PBx+PBx**2))*PQx*(QCx+QDx)*(15*(p+q)**2*f(3,t)+4*p*PQx**2*q*(5*(p+q)*f(4,t)+p*PQx**2*q*f(5,t)))/(p+q)**5+(4*PQx*q*(QCx+QDx)*(3+2*q*QCx*QDx)*(15*(p+q)**2*f(3,t)+4*p*PQx**2*q*(5*(p+q)*f(4,t)+p*PQx**2*q*f(5,t)))/(p*(p+q)**5)*(8*(PAx+PBx)*PQx*(3+q*(QCx**2+4*QCx*QDx+QDx**2))*(15*(p+q)**2*f(3,t)+4*p*PQx**2*q*(5*(p+q)*f(4,t)+p*PQx**2*q*f(5,t)))/(p+q)**5+(8*(PAx+PBx)*(QCx+QDx)*(15*(p+q)**3*f(3,t)+30*p*PQx**2*q*(p+q)*(3*(p+q)*f(4,t)+2*p*PQx**2*q*f(5,t))8*p**3*PQx**6*q**3*f(6,t)))/(p+q)**6+(2*(3+p*(PAx**2+4*PAx*PBx+PBx**2))*(15*(p+q)**3*f(3,t)+30*p*PQx**2*q*(p+q)*(3*(p+q)*f(4,t)+2*p*PQx**2*q*f(5,t))+8*p**3*PQx**6*q**3*f(6,t)))/(q*(p+q)**6)+(2*(3+q*(QCx**2+4*QCx*QDx+QDx**2))*(15*(p+q)**3*f(3,t)+30*p*PQx**2*q*(p+q)*(3*(p+q)*f(4,t)+2*p*PQx**2*q*f(5,t))+8*p**3*PQx**6*q**3*f(6,t)))/(p*(p+q)**6)$$

→ 787 MUL, 261 ADD, 69 FUNC

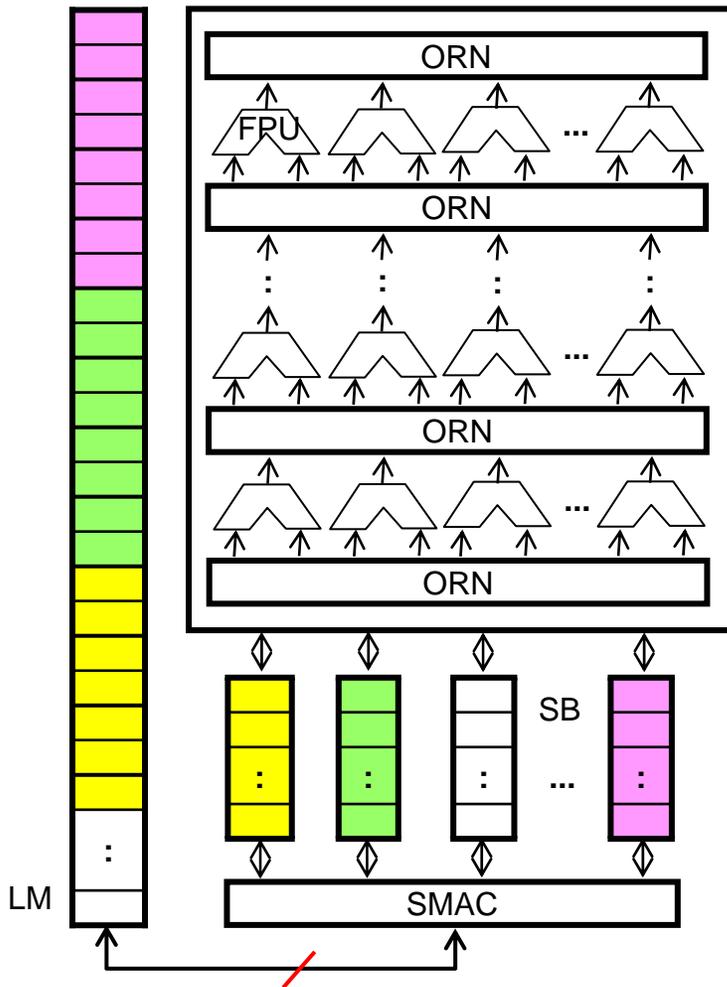
$$\text{tei}(3,1,1,1)=(PAy*(1+2*p*PAx*PBx)*(1+2*q*QCx*QDx)*f(0,t))/q+(((p+q)**4*((PAy+PQy)*q*(1+2*q*QCx*QDx)+p*(PAy+2*PAx*PAy*PQx*q+2*PAy*PBx*PQx*q+2*PAx*PBx*PQy*q+2*PAx*PAy*q*QCx+2*PAy*PBx*q*QCx+2*PAy*PQx*q*QCx+2*q*((PAy*(PAx+PBx+PQx))+2*(PAy*(PAx+PBx)*PQx+PAx*PBx*PQy)*q*QCx)*QDx)*2*p**2*PAx*PAy*PBx*(1+2*PQx*q*(QCx+QDx)))*f(1,t))/q+(p+q)*((p+q)*((p+q)*(3*p*PAy+6*p**2*PAx*PAy*PQx+6*p**2*PAy*PBx*PQx+2*p**2*PAy*PQx**2+4*p**3*PAx*PAy*PBx*PQx**2+p*PQy+2*p**2*PAx*PBx*PQy+2*p*PAy*PQx**2*q+PQy*q+2*p*PAx*PQx*PQy*q$$

→ 116 MUL, 31 ADD, 2 FUNC

RDPにマッピングするデータフローの例 ～(ps,ps)型積分計算の場合～

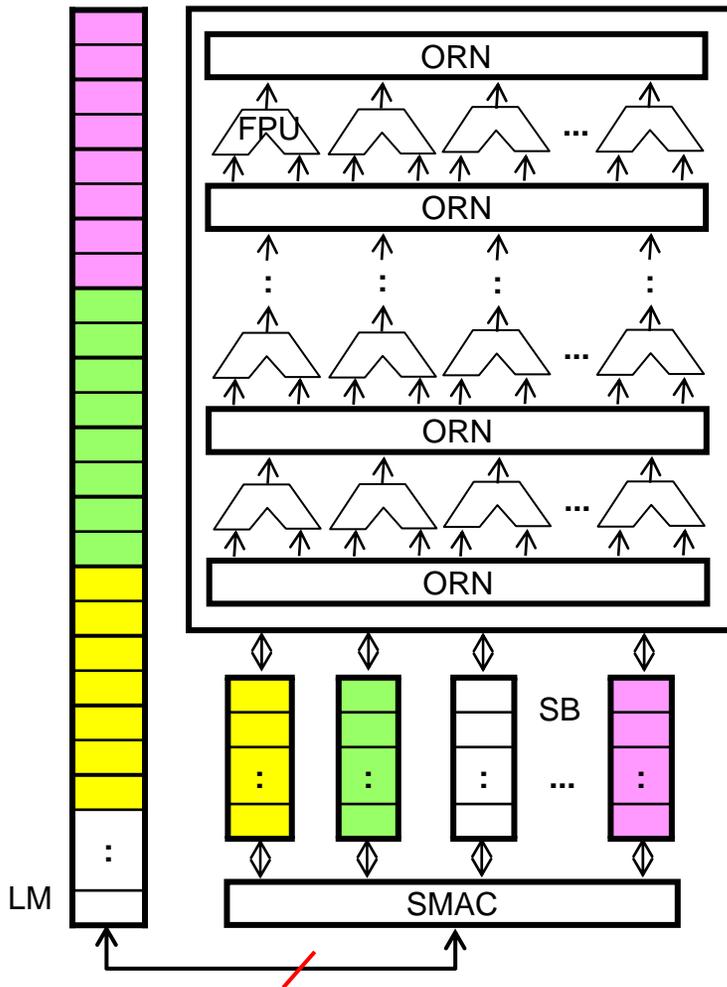


提案：再構成可能大規模データパス (RDP: Reconfigurable Datapath)



- 用途は？
 - 主プロセッサに対するコプロセッサ
 - コア計算部のループボディのデータフロー全体を直接マッピング
 - ORNでデータ依存関係を維持
 - 原理的には毎クロックサイクル、新しいイタレーションを実行開始可能
- 効能は？
 - 「必要とするメモリアクセス回数」を大幅に削減！
 - 従来プロセッサは、潜在的に「1演算につき3回のメモリアクセス」が必要

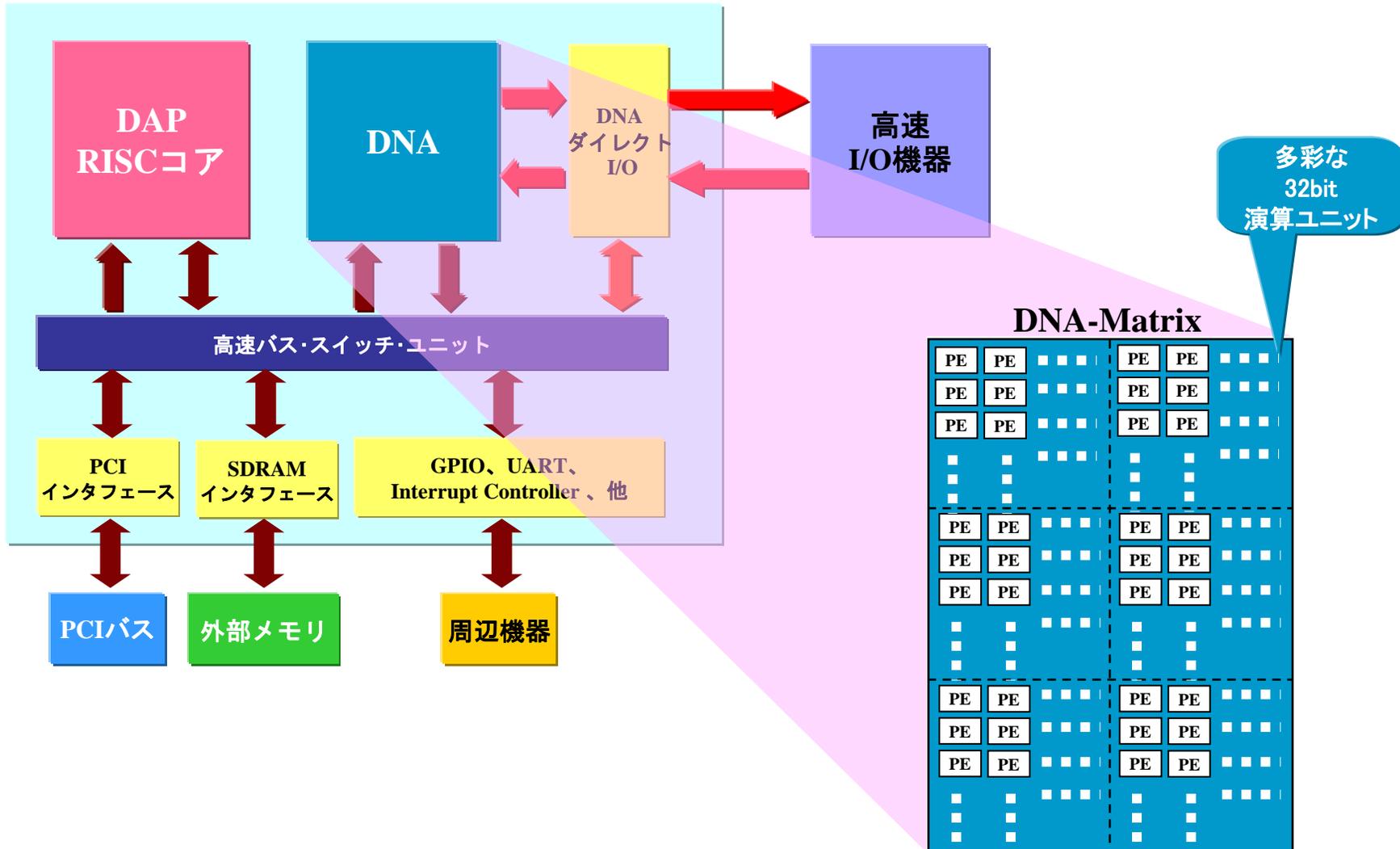
提案：再構成可能大規模データパス (RDP: Reconfigurable Datapath)



- 影響は？
 - プログラミングモデル：影響なし
 - コンパイラ最適化：
 - ループ・コラプシング (loop collapsing)：複数のループをまとめて、ループボディの大きなループを作る
- 継続性、経済的効果は？
 - 組み込みシステム、SoC業界との間での技術移転が可能
 - 例) IP FLEX DAP/DNA
 - 例) Stretch



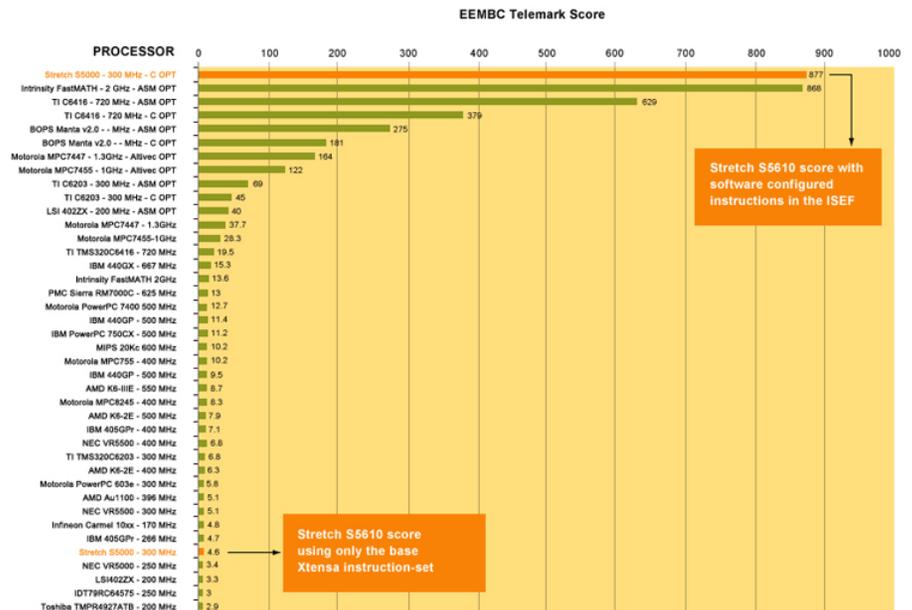
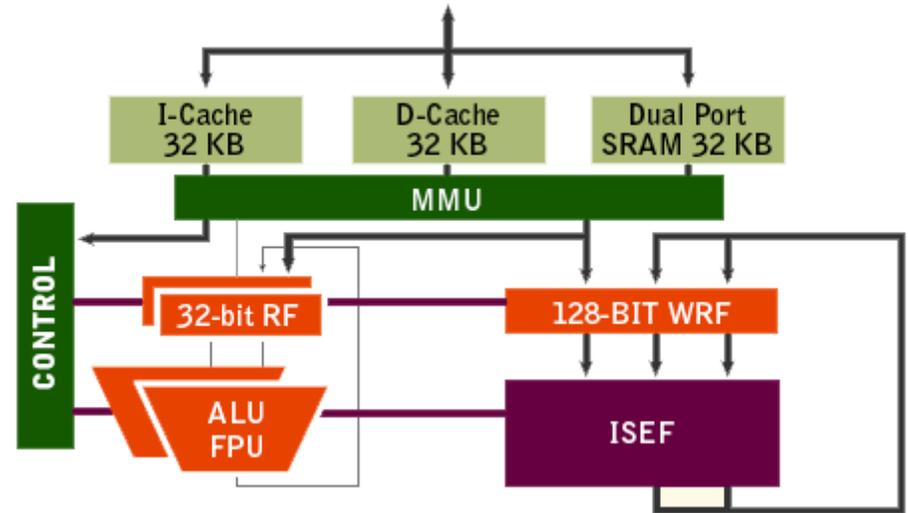
IP FLEX DAP/DNA2





Stretch

- 300 MHz, 32-bit Xtensa-based processor
- 16- and 24-bit instructions
- FPU
- MMU with TLB
- Stretch Instruction Set Extension Fabric
 - Aligned load and store
 - 8, 16, 32, 64, and 128 bit
 - Unaligned load and store
 - Up to 16 bytes variable byte streaming I/O
 - Up to 32 bits variable bit streaming I/O
- User-defined extensions to the core ISA
 - Defined in C/C++
 - Fully pipelined and interlocked
- Low power consumption
- Support for standard operating systems



10PFlop/sのマシンをどう作るか？

達成目標性能: LINPACK 10PFlop/s

↓ 仮定: 実行効率60%

ピーク性能: 16PFlop/s

仮定: クロック周波数1GHz

演算器数: 16M個

仮定: 4-way MUL&ADD SIMD

プロセッサコア数: 2M個

仮定: 4-way CMP

プロセッサチップ数: 512K個

プロセッサチップ性能: 32GFlop/s

仮定: 4B/s@Flop/s

要求オフチップメモリバンド幅: 128GB/s

仮定: クロック周波数1GHz

演算器数: 16M個

仮定: 1024-way RDP

プロセッサチップ数: 16K個

プロセッサチップ性能: 1TFlop/s

仮定: 0.0625B/s@Flop/s

要求オフチップメモリバンド幅: 64GB/s