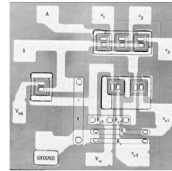
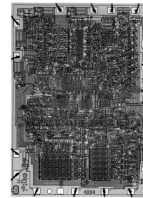


## Partie 1 - (Micro)processeurs

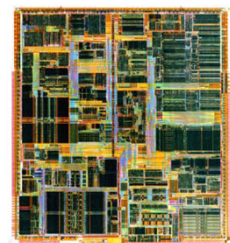
## Historique



Premier circuit intégré  
bipolaire, Motorola, porte 3  
entrées, 1966



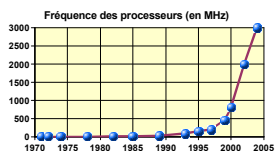
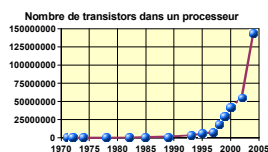
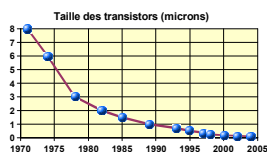
1971, Intel 4004,  
1000 transistors  
1 MHz NMOS-only tech  
pour le gain en vitesse



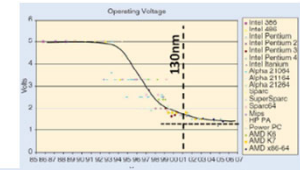
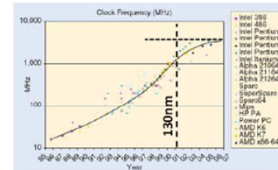
2000, Pentium 4 en CMOS (CMOS  
depuis 70 nm afin de réduire la puissance consommée)  
Multi-million portes

En numérique, le CMOS est aujourd'hui la technologie dominante

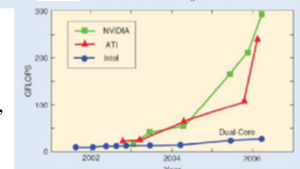
## "Loi" de Moore



## MAIS point d'inflexion !!!



L'évolution suivante  
est liée aux architectures  
(parallélisme niveau instruction ILP,  
multi-cœur ou spécifiques)



## Augmentation de la puissance de calcul

**Augmentation fréquence ?** Limites (conso., ...) => Réduction durée d'un traitement vs. durée d'un ensemble de traitements (séquentiel vs. parallèle), ou les deux ...

**Instruction-Level Parallelism (ILP)**

... + **organisation mémoire** (von Neumann vs. Harvard, caches ...)

**Thread-Level Parallelism (TLP)**

Dans un processeur,  
ou au niveau système  
("architectures parallèles")

**Data-Level Parallelism (DLP)**

**Co-processeurs ... logiciel vs. matériel**  
**FPGA et architectures reconfigurables**

## Puissance de calcul des microprocesseurs

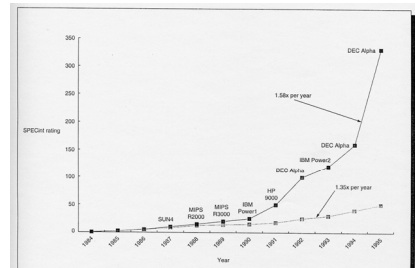
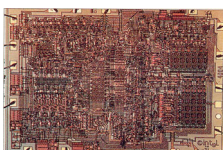


FIGURE 1.1 Growth in microprocessor performance since the mid 1980s has been substantially higher than in earlier years. This chart plots the performance as measured by the SPECint benchmarks. Prior to the mid 1980s, microprocessor performance growth was largely technology driven and averaged about 25% per year. The increase in growth since then is attributable to more advanced architectural ideas. By 1995 this growth leads to more than a factor of five difference in performance. Performance for floating-point-oriented calculations has increased even faster.

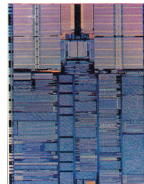
[Hennessy & Patterson 96]

## Du CPI à l'IPC



Nov. 1971 : introduction du 14004  
2300 transistors, 12mm<sup>2</sup>

25 ans



1996  
IBM présente le P2SC  
15M transistors (x6500), 335 mm<sup>2</sup> (x29),  
>x32 instructions/cycle, >x20.000 MIPS



Electronique  
International  
Hebdo  
31 Oct. 1996

## Processeurs et domaines d'application

Contraintes liées au domaine

Systèmes  
embarqués

PC/Serveur

- E/S multiples (standard ou spécialisées)
- Faible coût, sûreté, surface maximum, consommation, ...
- Puissance de calcul adaptée, efficacité énergétique (Perf/W)

- Puissance de calcul ...
- ... et autres

From ...

MCU  
8/16 bits ...

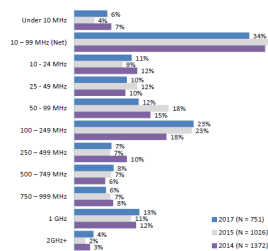
... to

Many-core CPU  
32/64 bits ...

MPU/CPU  
superscalaire / SMT  
32/64 bits ...

## Fréquences d'horloge

My current embedded project's main processor clock rate is:



The average processor clock rate was:  
445 MHz in 2017  
392 MHz in 2015  
420 MHz in 2014  
485 MHz in 2013

EE Times embedded

2017 Embedded Markets Study

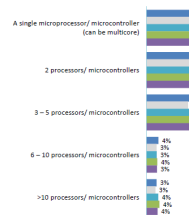
© 2017 Copyright by AsperCore. All rights reserved.

SEI - 2ème année

Architectures de processeurs et sécurité matérielle

## Nombre de processeurs

My current embedded project contains:



The average number microprocessor/microcontrollers per project was:  
2.3 in 2017  
2.1 in 2015  
2.4 in 2014  
2.4 in 2013  
2.3 in 2012

EE Times embedded

2017 Embedded Markets Study

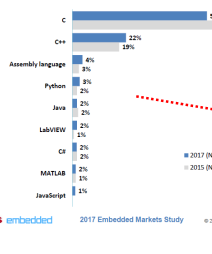
© 2017 Copyright by AsperCore. All rights reserved.

SEI - 2ème année

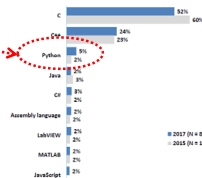
Architectures de processeurs et sécurité matérielle

## Langages de programmation

My current embedded project is programmed mostly in:



My next embedded project will likely be programmed mostly in:



EE Times embedded

2017 Embedded Markets Study

© 2017 Copyright by AsperCore. All rights reserved.

EE Times embedded

2017 Embedded Markets Study

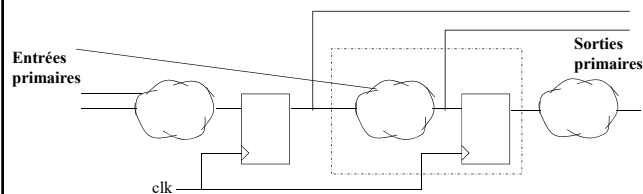
© 2017 Copyright by AsperCore. All rights reserved.

SEI - 2ème année

Architectures de processeurs et sécurité matérielle

## Circuit logique synchrone - rappel

- Hors mémoires : logique combinatoire entre deux "barrières" de registres + connexions à des E/S primaires
- Fréquence globale limitée par le bloc le plus "lent" entre deux barrières + synchronisation des entrées primaires



SEI - 2ème année

Architectures de processeurs et sécurité matérielle

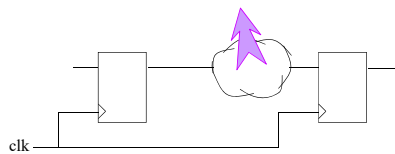
## Un sous bloc (ou étage) interne de base

- Logique combinatoire entre deux "barrières" de registres
- Comment augmenter la fréquence d'horloge ??

Exemple : calcul de la parité sur 16 bits

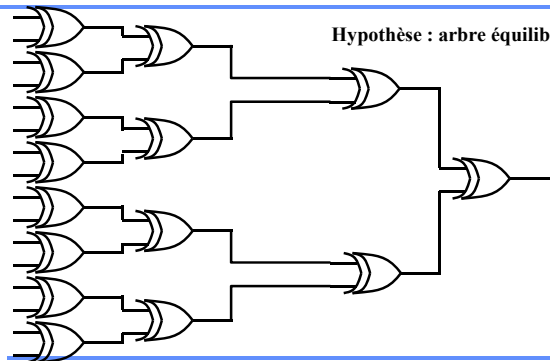
$$F = a \oplus b \oplus c \oplus d \oplus e \oplus f \oplus g \oplus h \oplus i \oplus j \oplus k \oplus l \oplus m \oplus n \oplus o \oplus p$$

Implantation sur une bibliothèque avec des portes XOR2

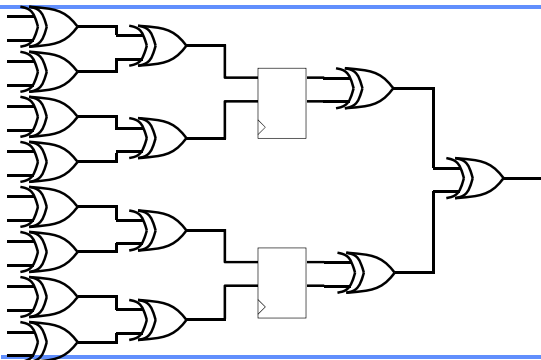


## Réalisation de la fonction de base

Hypothèse : arbre équilibré



## Ajout d'un pipeline 2 étages

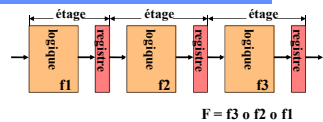


## Impact du pipeline

- 2 étages
  - ◆ Ajout de 4 bascules
  - ◆ Gain en fréquence ?
  - ◆ Gain en temps de calcul ?
  - ◆ Impact sur la latence ?
  - ◆ Gain en débit ?
- 4 étages ?
- Et 3 étages ?

## Structure pipeline - Résumé

### □ Partition + ajout de registres (composition de fonctions)



### □ Utile pour améliorer la fréquence maximum

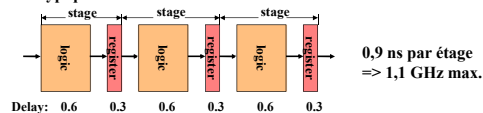
- ◆ Si la décomposition peut être équilibrée et si possibilité au niveau système (maillon faible = étage le plus lent, limitation supplémentaire potentielle par l'horloge système)
- ◆ Mais pas dans un rapport N avec N étages ( $T_{\text{setup}}/T_{H \rightarrow Q}$  des registres)

### □ N'améliore pas forcément les performances !

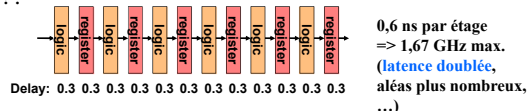
- ◆ Calculs multiples enchaînés possibles ? Sinon latence augmentée sans gain en puissance de calcul
- ◆ Taux d'utilisation des étages ( $\Rightarrow$  cf. aléas dans les processeurs)
- ◆ Energie (registres, arbres d'amplification), surface ...

## Impact des bascules : pipeline vs. superpipeline

### Micro-architecture typique du Pentium 3 :

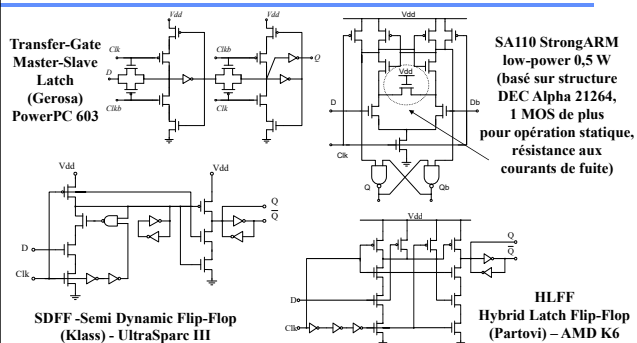


### Pentium 4 ? :



Augmentation de la fréquence  $\Rightarrow$  choix des éléments de mémorisation pour réduction des retards (compromis avec énergie, surface ...)

## Bascules D utilisées dans des processeurs



## Micro-parallélisme (ILP vers DLP/TLP)

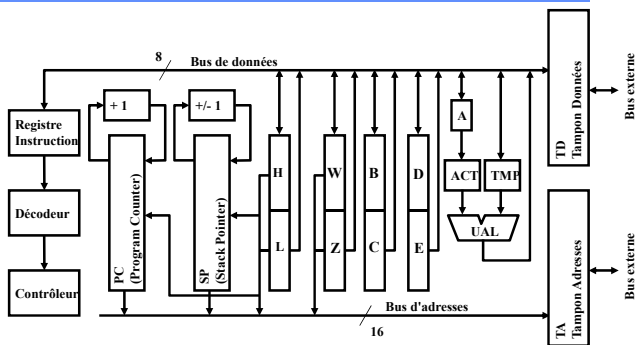
### □ Evolution progressive

- ◆ 8 bits : limitations technologiques  $\Rightarrow$  quasiment aucun parallélisme des opérations internes (exemple : 18080)
- ◆ Prefetch (exemple : MC68000)
- ◆ Queue d'instructions, séparation de l'unité d'exécution (exemple : 18088)
- ◆ RISC (ISA régulier, pipeline) + pipelines dans les CISC
- ◆ Superpipeline, superscalaire et VLIW
- ◆ Instructions SIMD (DLP)
- ◆ MT, SMT, ... (TLP)

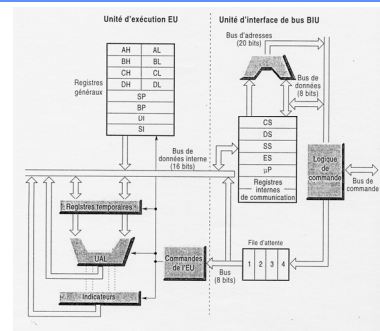
### □ Architecture Von Neumann vs. Harvard

- ◆ Séparation données/instructions
- ◆ Augmentation de la bande passante, parallélisation des accès mémoire
- ◆ Contraintes liées au nombre de broches

## Architecture du microprocesseur I8080



## Architecture du microprocesseur I8088



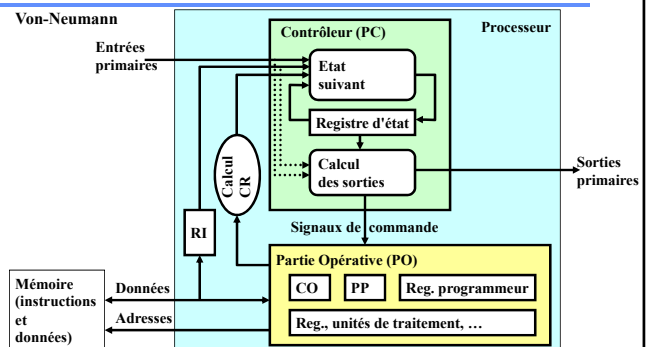
[Lilen 95]

## Modèles d'exécution

Modèle d'exécution	mémoire-mémoire	mémoire-accumulateur	mémoire-registre	registre-registre	pile
Génération (exemples)	ordinateurs (VAX 11)	µp 8 bits (I8080, M6800)	µp 16 bits (I8086, M68000)	µp RISC (I80860, M88000)	(Transputer)
Opérandes (mem. ext., spécifiées)	(3,3)	(1,1)	(1,2)	(0,3)	(0,0)
Inst. calcul	Mem op Mem -> Mem	ACC op Mem -> ACC	Rx op Mem -> Rx	Rx op Ry -> Rz	Pile1 op Pile2 -> Pile1
Exécution M1 op M2 -> M3	1 instruction Op M1, M2, M3	3 instructions Load M1, ACC Op M2 Store ACC, M3	3 instructions Load M1, Rx Op Rx, M2 Store Rx, M3	4 instructions Load M1, Rx Load M2, Ry Op Rx, Ry, Rz Store Rz, M3	4 instructions Push M1 Push M2 Op Pop M3
Cas particulier			(0,2): Rx op Ry -> Rx		

Note : le modèle registre-registre est aussi appelé "chargement-rangement" (seules les instructions Load et Store accèdent à la mémoire externe de données)

## Architecture de base d'un (micro-)processeur



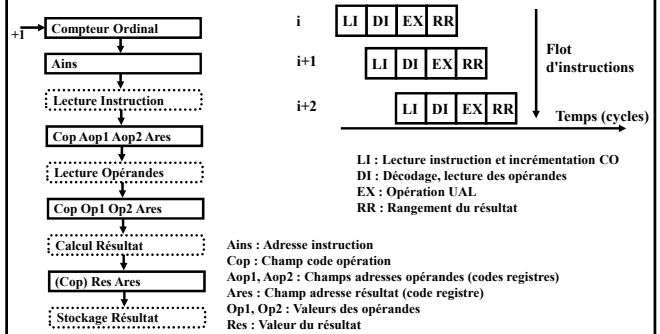
## Interprétation d'une (macro-)instruction

### Principales phases :

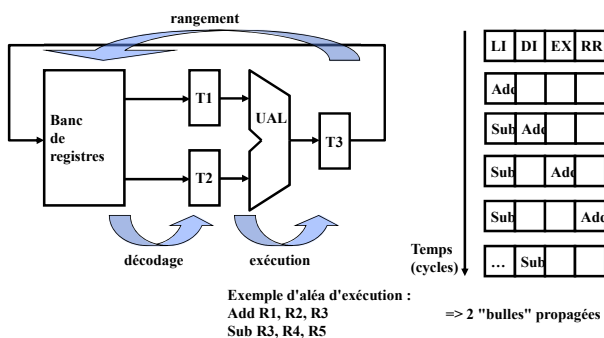
- ◆ Recherche du code de l'instruction (fetch), incrémentation CO
- ◆ Décodage
- ◆ Calcul d'adresses (selon mode d'adressage)
- ◆ Transfert des opérandes
- ◆ Exécution (traitement)
- ◆ Stockage du résultat

### Phases utiles dépendantes de l'instruction et du mode d'adressage (donc initialement du modèle d'exécution puis de la définition de l'ISA)

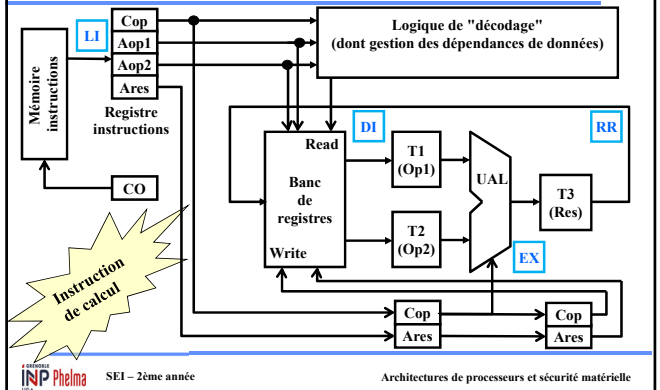
## Architecture pipeline : principe de base (Op calcul)



## Partie opérative pipelinée (principe)



## Structure pipeline pour calculs - PO + Contrôle



## Régularité du codage des instructions (ISA)

### Exemple : RISC V, instructions 32 bits "entier" de base

31	25	24	20	19	15	14	12	11	7	6	0	
funct7		rs2		rs1		funct3		rd		opcode		R-type <b>Registre-registre</b>
imm[11:0]				rs1		funct3		rd		opcode		I-type <b>Registre-immédiat</b>
imm[11:5]		rs2		rs1		funct3		imm[4:0]		opcode		S-type <b>Store, Branch</b>
				imm[31:12]				rd		opcode		U-type <b>Load imm., Jump</b>

The RISC-V ISA keeps the source (rs1 and rs2) and destination (rd) registers at the same position in all formats to simplify decoding. Except for the 5-bit immediates used in CSR instructions, immediates are always sign-extended, and are generally packed towards the leftmost available bits in the instruction and have been allocated to reduce hardware complexity. In particular, the sign bit for all immediates is always in bit 31 of the instruction to speed sign-extension circuitry.

<https://riscv.org/specifications/>

riscv-spec-v2.2



SEI – 2ème année

Architectures de processeurs et sécurité matérielle

## Aléas dans un pipeline (RISC ou CISC !)

### Aléas structurels

- ◆ Accès simultanés à une même ressource par plusieurs étages
- ◆ Exemple : accès à un cache unique données/instructions ...

### Aléas d'exécution

- ◆ Dépendances de données
- ◆ Branchements
- ◆ Exceptions (interruptions)

### Réduction du CPI apparent => nécessité de techniques de gestion des aléas pour limiter le nombre de "bulles"



SEI – 2ème année

Architectures de processeurs et sécurité matérielle

## Principales caractéristiques RISC

- **Modèle d'exécution chargement-rangement**
- **Format d'instruction de longueur fixe (1 mot) ou en tous cas régulier (pas "réduit" aujourd'hui !!)**,
- avec **modes d'adressage adaptés à la régularité du codage et à la limitation du nombre de bits pour coder une instruction => peu de modes d'adressage (surtout 2 : immédiat, base + déplacement M[Ri+dep])** avec possibilité de Load/Store par demi mots ou par octets, éventuellement avec décalages
- exemple MIPS : "LUI" Load Upper Immediate => chargement immédiat des 16 bits de poids fort d'un registre, avec mise à 0 des 16 bits de poids faible)
- "Grand" nombre de registres banalisés (32 ...)



SEI – 2ème année

Architectures de processeurs et sécurité matérielle

## Apparu RISC, fréquent dans les CISC récents

- "Grand" nombre de registres banalisés (32 ...)
- Pipeline interne / techniques de gestion des aléas
- Contrôleur (FSM) remplacé par les registres de contrôle (commandes internes) du pipeline
- Utilisation de caches pour réduire les temps d'accès mémoire
- CPI (apparent) de l'ordre de 1 ou inférieur

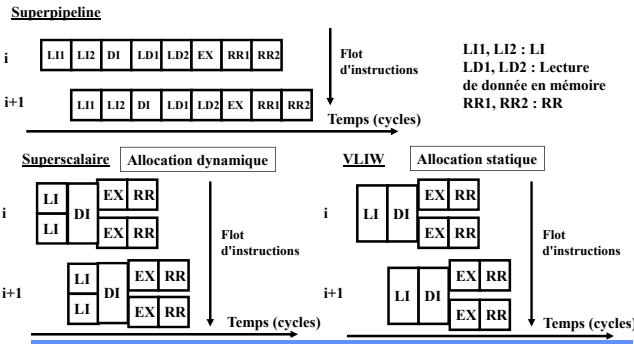


SEI – 2ème année

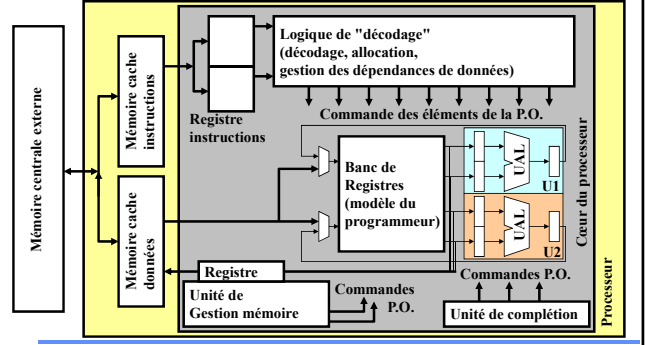
Architectures de processeurs et sécurité matérielle



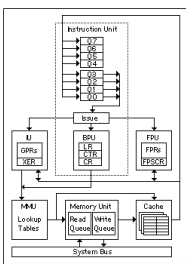
## Architecture pipeline : extensions (principes)



## Processeur superscalaire élémentaire

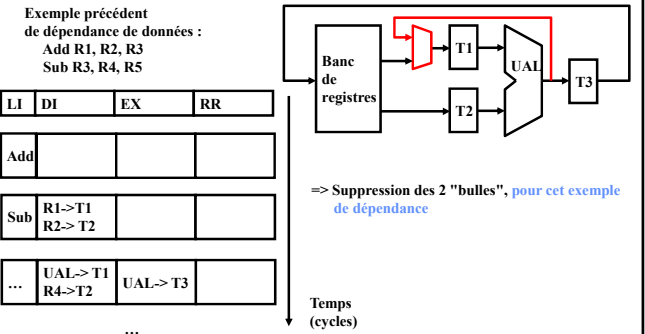


## PowerPC 601

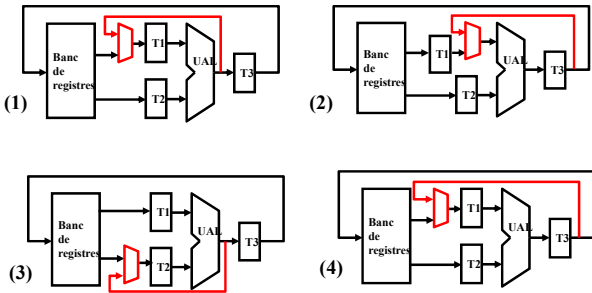


- ❑ Fetch : jusqu'à 8 instructions chargées en queue => allocation en désordre
- ❑ BPU : unité de prédiction de branchements, 2 étages de pipeline
- ❑ IU/FPU : traitements entiers/flottants, 4 et 6 étages de pipeline
- ❑ Accès mémoire (Load/Store) : 5 étages de pipeline
- ❑ Cache : associatif 8 voies, unifié instructions/données
- ❑ MMU : adressage jusqu'à 4 Gb de mémoire physique et 4 Peta-octets de mémoire virtuelle (2<sup>52</sup> octets ... ou 6,2 millions de CD-ROM !)

## Traitement d'une dépendance de données



## Modifications équivalentes ??



## Traitement d'un aléa d'exécution (branchement)

### Trois grandes approches

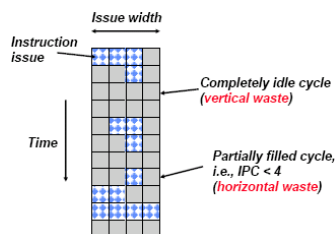
- ◆ **Branchements retardés (jeu d'instructions des premiers RISC)** : exécution de l'instruction suivant le Br quelque soit la valeur de la condition – "Delay slot"
- ◆ **Anticipation de branchement (ex. PowerPC 601)** : si la condition est connue, remplacement de l'instruction Br par l'instruction destination du saut, dans la file d'attente d'instructions, avant le décodage effectif ("issue") => pas de pénalité dans le pipeline, et la lecture d'instructions continue en séquence avec l'instruction destination. Possibilité aussi de 2 queues d'instructions avec sélection lorsque la condition est connue (SuperSparc).
- ◆ **Prédiction de branchement : exécution spéculative.**
  - Prédiction statique : règle prédéfinie
  - Prédiction dynamique : mémorisation des adresses des instructions de branchement exécutées et des adresses de destination calculées pour chacune d'elles

- **Remarque : possibilité en compilation de "dérouler" les boucles avec condition d'arrêt fixe (boucles for) pour supprimer le branchement ... Ou utilisation de "boucles matérielles"**

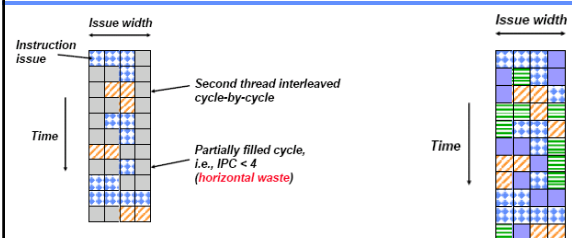
## Parallélisme niveau tâches : motivation

- **Exécution d'une tâche unique : utilisation limitée des ressources disponibles en fonction**

- ◆ Des types d'instructions à exécuter par rapport aux unités disponibles
- ◆ Des dépendances et autres aléas (dont accès mémoire)



## Parallélisme niveau tâches : gains

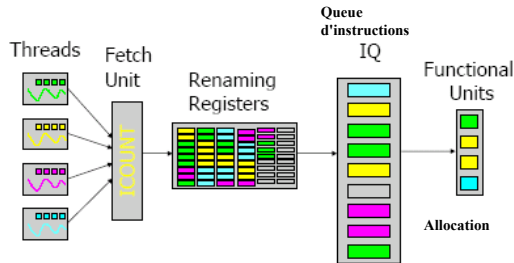


MT ("vertical")  
=> Gain sur les cycles perdus

SMT "idéal"  
=> Utilisation maximisée

Schémas issus de [www.cs.ucr.edu/~bhuyan/cs162]

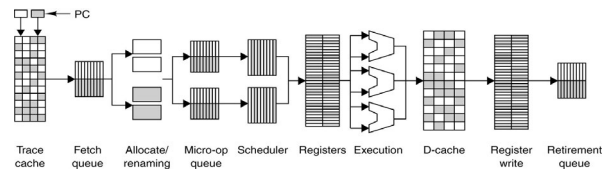
## Du superscalaire au processeur SMT



Renommage de registres : utilisé en superscalaire pour supprimer les fausses dépendances, exploité aussi pour les fausses dépendances entre tâches (utilisation de registres supplémentaires, non visibles dans l'ISA)

## Microarchitecture Pentium 4 NetBurst

- Partage de ressources (unités de traitement, gestion des aléas, mémoires ...)
- Augmentation du nombre d'étages de pipeline mais différent du superpipeline (sous-tâches distinctes entre étages)



## Parallélisme niveau tâches : coûts

- Gestion de l'état de chaque tâche et maintien de la cohérence globale
  - ◆ Compteur de programme par tâche
  - ◆ Pile (et pointeur de pile) par tâche
  - ◆ Gestion des identificateurs de tâche sur l'ensemble des étapes de traitement (incluant les tables d'allocation des caches)
  - ◆ Gestion des queues d'instructions par tâche
  - ◆ Gestion des complétions par tâche
- Pertes de performances dues aux changements de contexte

## Parallélisme niveau tâches - vue d'ensemble

- MT : multi-threading
  - ◆ "Entrelacement" des recherches d'instructions sur plusieurs tâches
  - ◆ Gestion en parallèle de plusieurs environnements d'exécution
  - ◆ Changement possible de contexte entre tâches : grain fin (chaque cycle), sur un aléas d'accès au cache niveau 1, sur un aléas d'accès au cache niveau 2, sur attente d'une requête externe ...
- SMT : simultaneous multi-threading
  - ◆ Extension du MT => chargement simultané dans la queue d'instructions, à chaque cycle, d'instructions provenant de plusieurs process ou tâches
  - ◆ Permet d'augmenter l'utilisation des différentes ressources de traitement (ex. : unités de calcul entier et de calcul flottant)
- HT : hyper-threading technology (Intel) => Pentium 4 à 3,06 GHz
  - ◆ SMT + techniques de prédiction au niveau tâche
  - ◆ + techniques de parallélisation des accès dans la hiérarchie mémoire (MLP - Memory Level Parallelism)

=> Lien très étroit avec les systèmes d'exploitation !!!!

# Comparaison avec multi-cœur (débit)

The diagram compares four multi-processor architectures using a 5x5 grid of 'Issue slots' to represent execution over time (vertical axis) and processors (horizontal axis).

- Superscalar horizontal waste:** Shows a single thread (blue) executing across multiple processors. Horizontal waste is indicated by empty slots in a row, and vertical waste is indicated by empty slots in a column.
- Traditional Multithreading:** Shows multiple threads (blue, red, green) interleaved across processors. Horizontal waste is indicated by empty slots in a row.
- Single-chip Multiprocessor:** Shows multiple threads (blue, red, green) interleaved across processors. Horizontal waste is indicated by empty slots in a row.
- SMT:** Shows multiple threads (blue, red, green, magenta, cyan) interleaved across processors. Horizontal waste is indicated by empty slots in a row.

Legend:

- Thread 1 (Blue)
- Thread 2 (Red)
- Thread 3 (Green)
- Thread 4 (Magenta)
- Thread 5 (Cyan)

Schéma issu de [www.cs.ucr.edu/~bhuyan/cs162]

INP Phelma SEI - 2ème année Architectures de processeurs et sécurité matérielle

# Parallélisme niveau tâches : exemples

- MIPS32 – Famille 34K
  - ◆ Processeurs uniscalaires (1 résultat max / cycle), 9 niveaux de pipeline
  - ◆ Mais SMT (au niveau de l'étage de lecture), jusqu'à 5 tâches (micro-tâches, ou programmes distincts) traitées simultanément dans le pipeline
  - ◆ Réduction du coût des changements de contexte, possibilité de masquer les attentes dues aux opérations lentes (ex. accès mémoire données)
  - ◆ Conserve la compatibilité avec les logiciels à tâche unique ("single-threaded")
- Sun : processeur Niagara 2, destiné aux serveurs
  - ◆ 64 tâches traitées simultanément ...
  - ◆ Possibilité de coupler 2 processeurs pour traiter 128 tâches simultanément de manière cohérente.
  - ◆ Problème de gestion par l'OS, contraintes de synchronisations entre tâches, limitation de la bande passante interne et externe => quelle est le niveau limite d'efficacité du parallélisme niveau tâche ??

[Microprocessor Report, février 2006]

à l'origine de  
**INP Phelma**  
UNIVERSITÉ

SEI – 2ème année

Architectures de processeurs et sécurité matérielle

# Caractéristiques de processeurs du commerce

---

- Exemples (de 1993 ...)
  - ◆ Pentium  
32 bits, CISC superscalaire, 3,1 MTransistors, 294 mm<sup>2</sup> CMOS 0,8μ  
60-66 MHz – Conso 12-16W – SpecInt92 64,5 – SpecFp92 56,9  
\$900 par 1000 pièces
  - ◆ PowerPC 601  
32 bits, RISC superscalaire, 2,8 MTransistors, 120 mm<sup>2</sup> CMOS 0,65μ  
66-80 MHz – Conso 7W (à 66 MHz) – SpecInt92 80 – SpecFp92 105  
\$550 par 1000 pièces

## Où est l'erreur ?

---

COMMISSARIAT GÉNÉRAL  
**INP** Phelma  
1994

SEI – 2ème année

Architectures de processeurs et sécurité matérielle

# Evolution de l'architecture ARM

**Architecture RISC :**

- jeu d'instructions de longueur fixe: 32 bits
- architecture 'load/store'
- banc de registres d'usage général

Timeline of ARM architectures:

- V4 (Blue):** ARM7TDMI (1994), StrongARM (1996), ARM720T (1998), ARM920T (2000).
- ARMv5 (Green):** ARM1020 (2002), ARM1022E (2004), ARM1026J (2004), ARM1110J (2006).
- ARMv6 (Yellow):** XScale™ (2002), V6 cores (2006).

**Quelques spécificités :**

- jeu d'instructions Thumb sur 16 bits (puis Thumb2)
- technologie Jazelle : multiplie la vitesse d'exécution de code Java par 8
- extensions DSP (instructions supplémentaires)

ARM Limited INP Phelma SEI – 2ème année

Architectures de processeurs et sécurité matérielle

## Architectures de processeurs et sécurité matérielle

## Architectures de processeurs et sécurité matérielle

## Architectures de processeurs et sécurité matérielle

## Architectures de processeurs et sécurité matérielle

