

# Basse consommation dans les systèmes embarqués

Daniel CHILLET

ENSSAT - Université de Rennes 1  
INRIA/IRISA - Projet R2D2  
6 Rue de Kérampont - F-22300 LANNION

chillet@enssat.fr



## Plan

- ⌘ Evolution technologique
- ⌘ Pourquoi y-a-t il consommation ?
- ⌘ Puissance statique & dynamique
- ⌘ Qu'est ce qui consomme ?
- ⌘ Solutions :
  - ☒ au niveau technologique
  - ☒ au niveau logique
  - ☒ au niveau architectural
  - ☒ au niveau logiciel (non traitées ici)
  - ☒ au niveau système (non traitées ici)
- ⌘ Conclusions



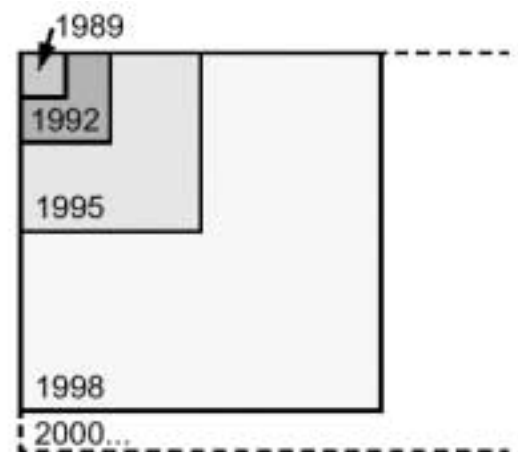
Année d'introduction	1997	1999	2001	2003	2006	2009	2012
Technologie (µm)	0,25	0,18	0,15	0,13	0,1	0,07	0,05
Fréquence (MHz)	750	1250	1500	2100	3500	6000	10000
Mémoire		256M	1G	4G	16G	64G	256G
Nombre de transistors	11M	21M	38M	77M	202M	520M	1350M
Coût par M Tr (Cent.)	500	290	166	97	42	18	8

## ⌘ Scaling technologique à chaque génération

- ☒ Fréquence augmente de 43%
- ☒ Capacité totale et tension d'alimentation sont réduites de 30%
- ☒ Énergie réduite de 65%
  - ☒  $E = C \cdot V_{dd}^2 = C \cdot 0.7 \cdot (V_{dd} \cdot 0.7)^2 = 0.35 \cdot C \cdot V_{dd}^2 = 35\% E'$
- ☒ Puissance réduite de 50%
  - ☒  $P = f \cdot C \cdot V_{dd}^2 = 1.43 \cdot f \cdot 0.35 \cdot C \cdot V_{dd}^2 = 50\% P'$
  - ☒ On considère que l'activité du circuit est constante

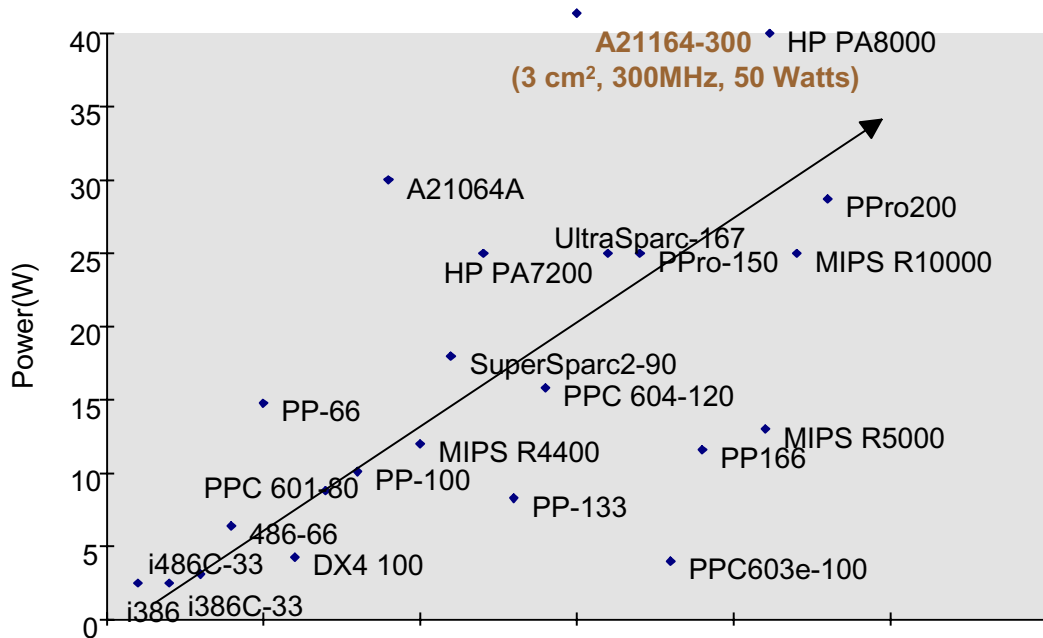
**Mais à nombre de transistors constant!**

- ☒ Densité de transistor double à chaque génération
- ☒ Surface des puces augmente de 25%
- ☒ Densité de puissance augmente avec un facteur 2
- ☒ Courant d'alimentation augmente de façon importante



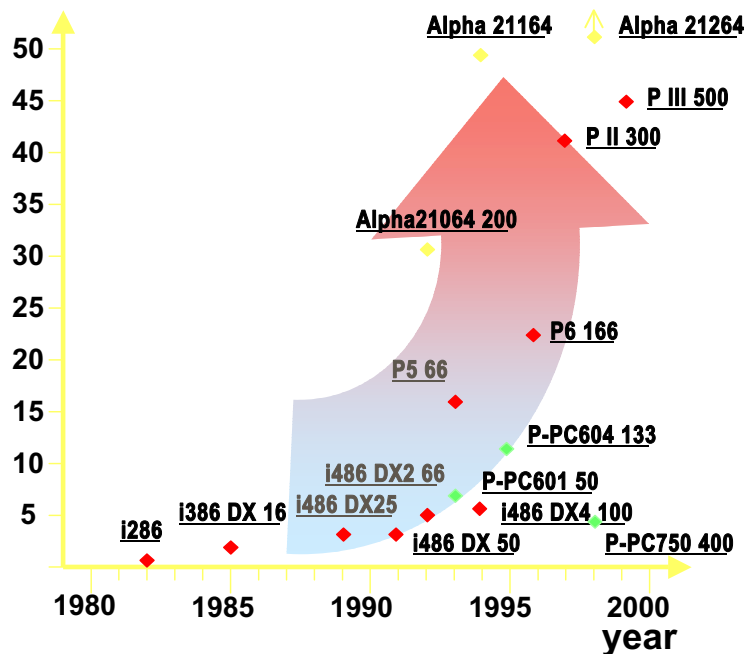
# Effet sur la consommation

⏏ x4 tous les 3 ans



# Evolution de la consommation

Power(W)



# Evolution de la consommation



⌘ e.g. Famille DEC/Compaq [Herrick99]

## ⏏ EV4

- ⊗ 200 MHz @ 3.3V
- ⊗ 30W @ 200 MHz
- ⊗ 1,7 Million Transistors
- ⊗ 233 mm<sup>2</sup>

## ⏏ EV5 (21164)

- ⊗ 350 MHz @ 3.3V
- ⊗ 60W @ 350 MHz
- ⊗ 9,3 Million Transistors
- ⊗ 298 mm<sup>2</sup>

## ⏏ EV6 (21264)

- ⊗ 575 MHz @ 2.2V
- ⊗ 90W @ 575 MHz
- ⊗ 15,2 Million Transistors
- ⊗ 314 mm<sup>2</sup>

## ⏏ EV7 (21364)

- ⊗ > 1000 MHz @ 1.5V
- ⊗ 100W @ 1GHz
- ⊗ ~100 Million Transistors
- ⊗ ~350 mm<sup>2</sup>

## ⏏ EV8

- ⊗ > 1-2 GHz (0.125 micron)
- ⊗ <150W
- ⊗ ~250 Million Transistors

# Effets sur les circuits



⌘ Dissipation chaleur :

- ⏏ diminution de la fiabilité  
MTBF/2 pour une augmentation de 10°C

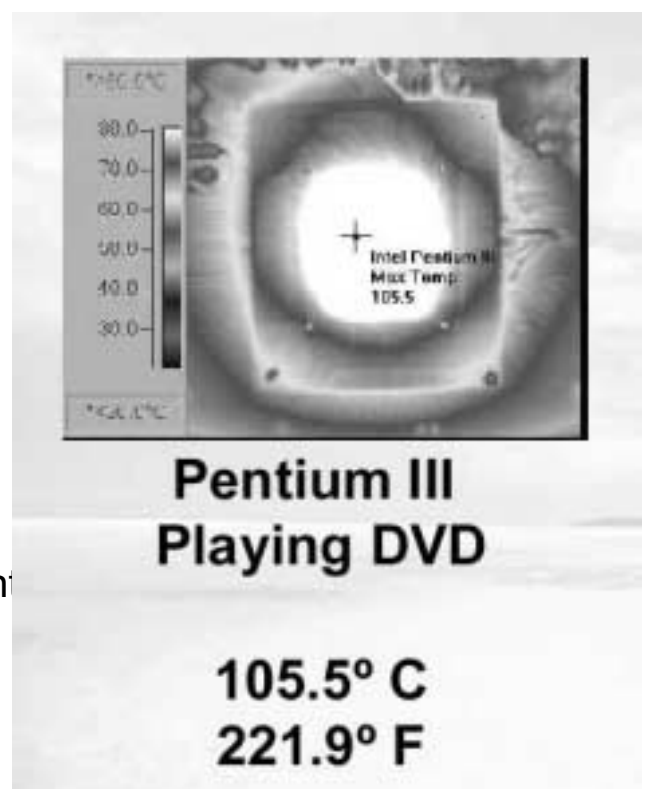
⌘ Augmentation du coût de la mise en œuvre

- ⏏ 1€/W si >40W

⌘ Augmentation volume

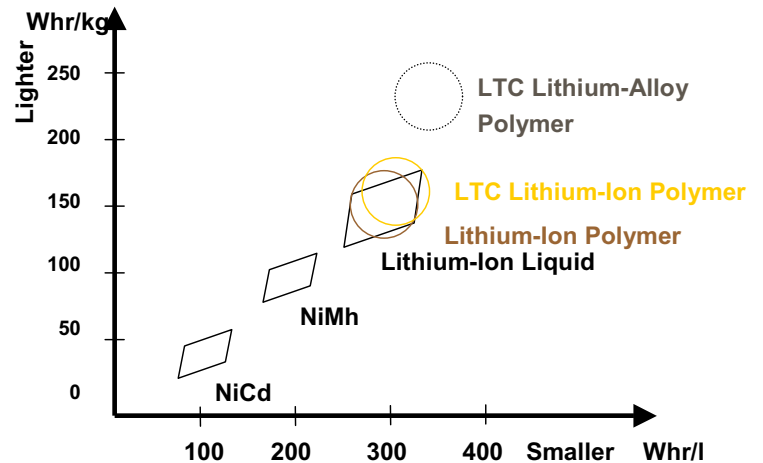
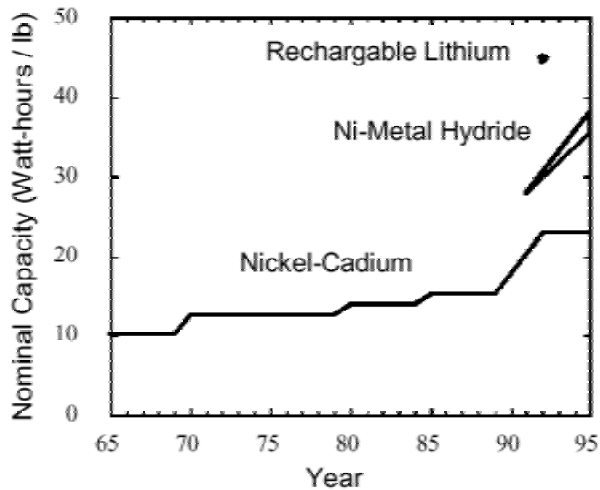
- ⏏ radiateur, ventilateur (volume sonore)

⌘ Diminution du temps de fonctionnement  
ou  
augmentation du poids des batteries



# Evolution de la technologie des batteries

## Evolution techno des batteries



[P. Senn 2000]

# Evolution de la technologie des batteries

Techno	NiCd Nickel Cadium	NiMh Nickel Metal Hybride	Li-ion Lithium Ion	Li-poly
Apparition	1956	1990	1992	1996
Tension délivrée (V)	1,2	1,2	3,6	3,7
Epaisseur (mm)	> 6	>6	>6	3
Capacité (Whr/Kg)	30-50	60-90	70-140	115-140
Durée de vie (cycles)	~1000	~ 1000	~ 500	~ 500
Extralopation pour 1kg à 20 W	22 - 27 h	40 - 45 h	50 - 55 h	65 - 70 h



## ⌘ Systèmes :

### ☑ capacité de traitement élevée :

- ☒ multimédia (visioconf, Mpeg, ...)
- ☒ interface conviviale (reco vocale, stylo, synthèse vocale, ...)
- ☒ sécurité des info (cryptographie, authentification, ...)

### ☑ portables

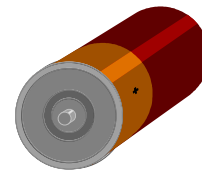
### ☑ faible coût

### ☑ peu volumineux

### ☑ peu lourd

### ☑ peu bruyant

### ☑ longue durée d'utilisation sans rechargement des batteries



# Pourquoi il y a consommation ?

## ⌘ Puissance = puissance statique + puissance dynamique

### ⌘ Puissance statique :

- ☑ même bloqué, les transistors ne sont pas complètement fermés

### ⌘ Puissance dynamique :

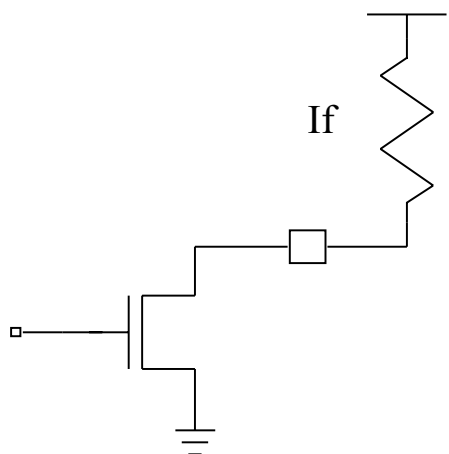
- ☑ lors de la commutation des transistors, il va se créer des court-circuits entre Vdd et la masse
- ☑ charge et décharge de capacités

# Puissance statique

## ⌘ Courants de fuite (Subthreshold Leakage Current)

- ☒ Même lorsque  $V_{gs} < V_t$  les transistors MOS ne sont pas équivalents à un interrupteur fermé
- ☒ Si on diminue  $V_{dd}$  ( $\neq V_t$ ) alors ces courants peuvent devenir significatifs

## ⌘ Courants inverses de diode Source/Drain - Substrat

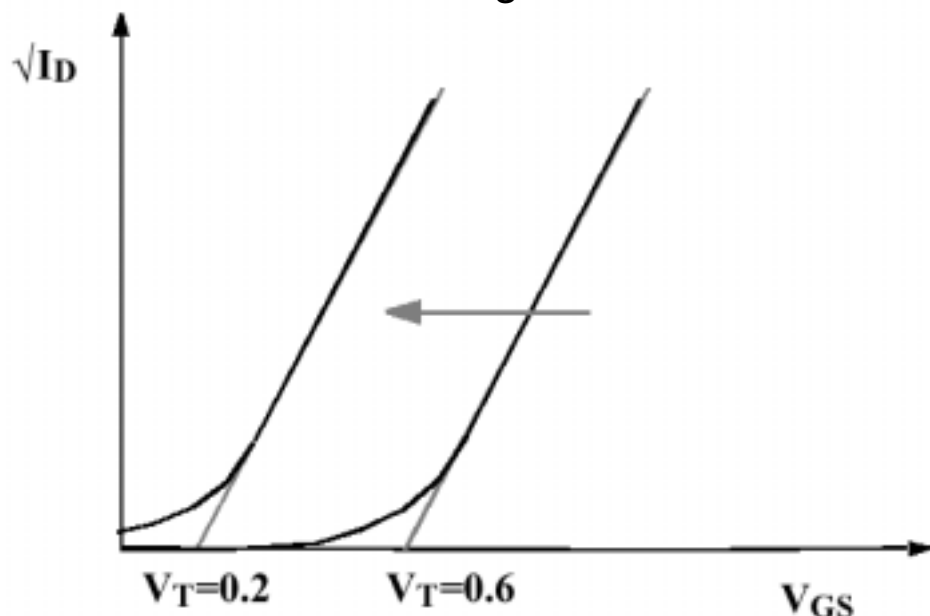


➡ Négligeables tant que  $V_{dd} > V_t$

# Puissance statique

## ⌘ Impact de la tension de seuil :

- ☒ diminution tension de seuil → augmentation du courant de fuite



- ☒ donc, pour un design low power, il vaut mieux une tension de seuil élevé

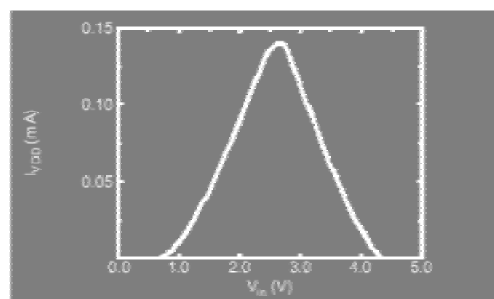
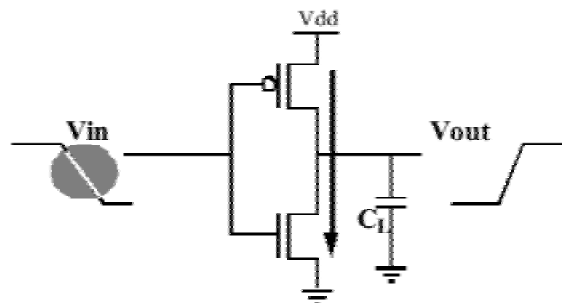
## ⌘ Négligeable ou pas ?

- ☒ souvent négligée pour les blocs logiques ayant beaucoup d'activité
- ☒ mémoires de grandes tailles, peu d'activité (relativement au nombre de transistor) → puissance non négligeable
- ☒ techno sub-micronique profonde :
  - ☒ augmentation de l'influence des interconnexions :
    - charge supplémentaire pour les portes logiques

# Puissance dynamique

## ⌘ Apparition lors de la commutation des transistors :

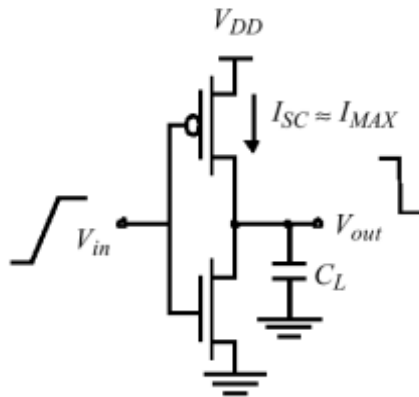
- ☒ Courants de court-circuit :  $I_{cc}$ 
  - ☒ Chemin de court circuit pendant la commutation des structures logiques
  - ☒ Pour les circuits *conçus correctement*  $\approx 15\%$



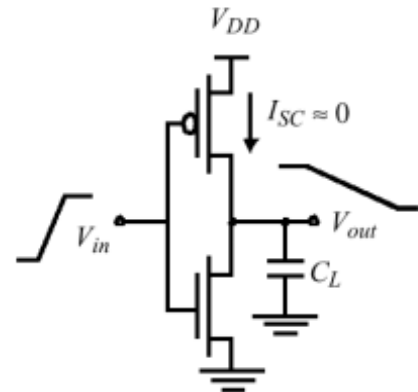


⚡ Courants de court-circuit : impact des fronts lents :

- ⊗ temps de court circuit plus long
- ⊗ dépend de la charge de la porte

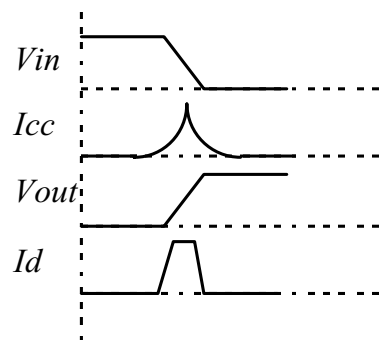
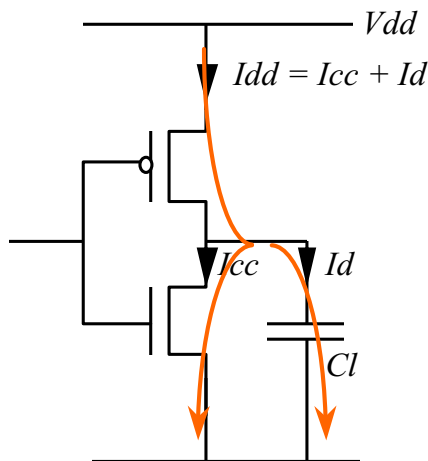


Capacité faible



Capacité élevée

⚡ Courants de charge/décharge des capacités :  $I_d$



# Puissance dynamique

- ⌘ Energie par transition =  $C_l \cdot V_{dd}^2$
- ⌘ Puissance = Energie par transition • Vitesse de transition  
=  $C_l \cdot V_{dd}^2 \cdot f$
- ⌘ La puissance n'est pas fonction de la **taille des transistors** de la cellule considérée

$$\begin{aligned} P_{dyn} &= C_l \cdot V_{dd}^2 \cdot f_{0 \rightarrow 1} \\ &= C_l \cdot V_{dd}^2 \cdot f \cdot P_{0 \rightarrow 1} \\ &= \alpha \cdot C_l \cdot V_{dd}^2 \cdot f \\ &= C_{EFF} \cdot V_{dd}^2 \cdot f \end{aligned}$$

- $C_{EFF}$  est la capacité effective =  $C_l \cdot P_{0 \rightarrow 1}$



# Puissance dynamique

- ⌘ e.g. porte inverseuse :

A	OUT
0	1
1	0

Fonction NOT

$$\begin{aligned} P(0 \rightarrow 1) &= P(\text{OUT} = 0) \cdot P(\text{OUT} = 1) \\ &= 1/2 \cdot 1/2 = 1/4 \end{aligned}$$

$$\alpha C = C_{eff} = 1/4 \cdot C_l$$

# Puissance dynamique

⌘ e.g. porte NOR statique à 2 entrées

A	B	OUT
0	0	1
0	1	0
1	0	0
1	1	0

Fonction NOR

Hypothèse :  $P(A=1) = 1/2$  et  $P(B=1) = 1/2$

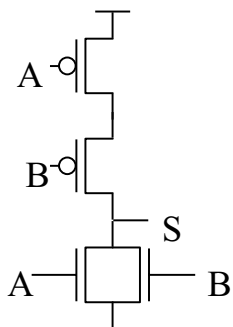
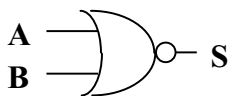
$$\begin{aligned}
 P(\text{OUT}=1) &= 1/4 \\
 P(0 \rightarrow 1) &= P(\text{OUT} = 0) \cdot P(\text{OUT} = 1) \\
 &= 3/4 \cdot 1/4 = 3/16
 \end{aligned}$$

$$\alpha C = C_{\text{eff}} = 3/16 \cdot C_I$$

# Puissance dynamique

⌘ Probabilités de transition d'une porte NOR

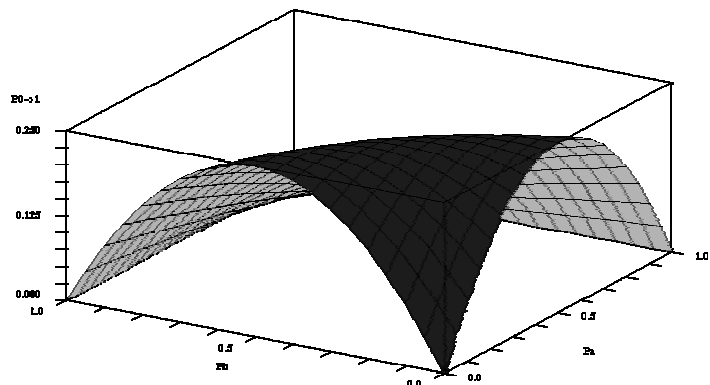
A	B	S
0	0	1
0	1	0
1	0	0
1	1	0



$$P_S = (1 - P_A)(1 - P_B) \quad \text{avec } P_A = P(A=1) \text{ et } P_B = P(B=1)$$

$$P_1 = P(S=1) = (1 - P_A)(1 - P_B)$$

$$\begin{aligned}
 P_{0 \rightarrow 1} &= P_0 \cdot P_1 = P(S=0) \cdot P(S=1) = (1 - P(S=1)) \cdot P(S=1) \\
 &= (1 - (1 - P_A)(1 - P_B)) (1 - P_A)(1 - P_B)
 \end{aligned}$$



⌘  $\alpha$  est une fonction dépendant fortement de la statistique des signaux

## ⌘ Probabilités de transition d'une porte XOR

$$P_S = P_A + P_B - 2 \cdot P_A \cdot P_B \quad \text{avec } P_A = P(A=1) \\ \text{et } P_B = P(B=1)$$

$$P_1 = P(S=1) = (1 - P_A)(1 - P_B)$$

$$P_{0 \rightarrow 1} = P_0 \cdot P_1 = P(S=0) \cdot P(S=1) = (1 - P(S=1)) \cdot P(S=1) \\ = (1 - (P_A + P_B - 2 \cdot P_A \cdot P_B)) \cdot (P_A + P_B - 2 \cdot P_A \cdot P_B) \\ = (1 - (\frac{1}{2} + \frac{1}{2} - 2 \cdot \frac{1}{2} \cdot \frac{1}{2})) (\frac{1}{2} + \frac{1}{2} - 2 \cdot \frac{1}{2} \cdot \frac{1}{2}) \\ = (1 - (1 - \frac{1}{2})) (1 - \frac{1}{2}) = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$$

A	B	S
0	0	0
0	1	1
1	0	1
1	1	0

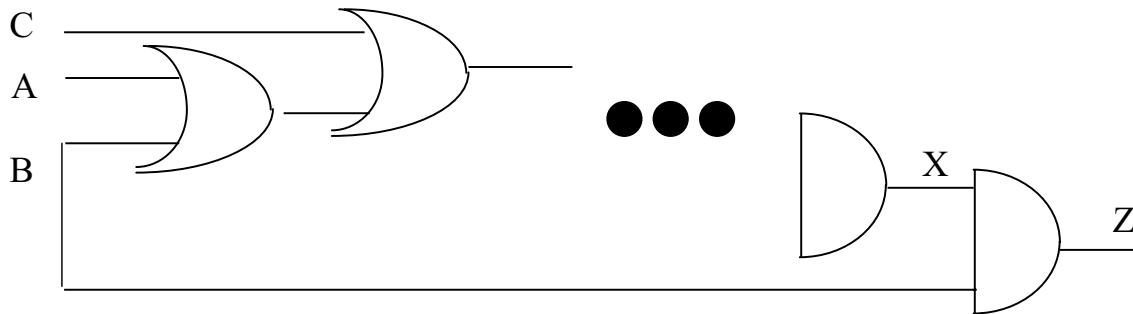
## ⌘ Fonction de la statistique des entrées

	$P_{0 \rightarrow 1}$	Ceff
AND	$(1 - P_A P_B) P_A P_B$	3/16
OR	$(1 - P_A)(1 - P_B) (1 - (1 - P_A)(1 - P_B))$	3/16
XOR	$(1 - (P_A + P_B - 2 P_A P_B)) (P_A + P_B - 2 P_A P_B)$	1/4
NOR	$(1 - (1 - P_A)(1 - P_B)) (1 - P_A)(1 - P_B)$	3/16

# Puissance dynamique

## ⌘ Problème de la reconvergence

☑ devient très rapidement complexe



$$P(Z=1) = P(B=1) \cdot P(X=1 \mid B=1)$$

# Puissance dynamique

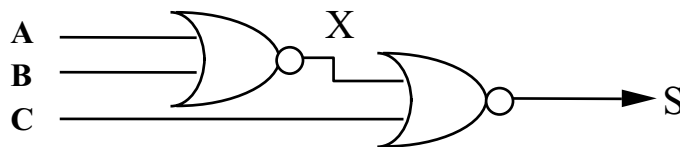
## ⌘ Problème des glitches

☑ non pris en compte dans les probabilités

$$P(A) = \frac{1}{2}$$

$$P(B) = \frac{1}{2}$$

$$P(C) = \frac{1}{2}$$



$$P(X=1) = 1/4$$

$$P(S=1) = 1/2 * 3/4 = 3/8$$

$$\alpha X = P(X=0) \cdot P(X=1)$$

$$= (1-P(X=1)) \cdot P(X=1)$$

$$= (1 - 1/4) \cdot 1/4$$

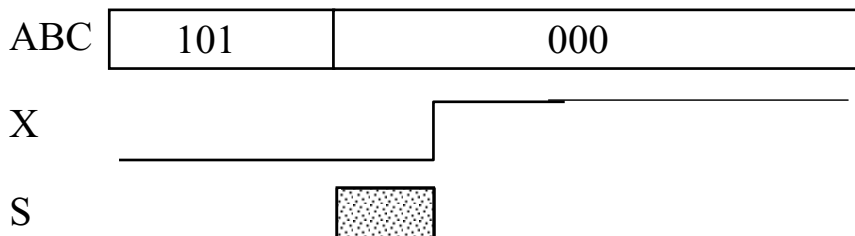
$$= 3/4$$

$$\alpha S = P(S=0) \cdot P(S=1)$$

$$= (1 - P(S=1)) \cdot P(S=1)$$

$$= (1 - 3/8) \cdot 3/8 = 5/8 \cdot 3/8$$

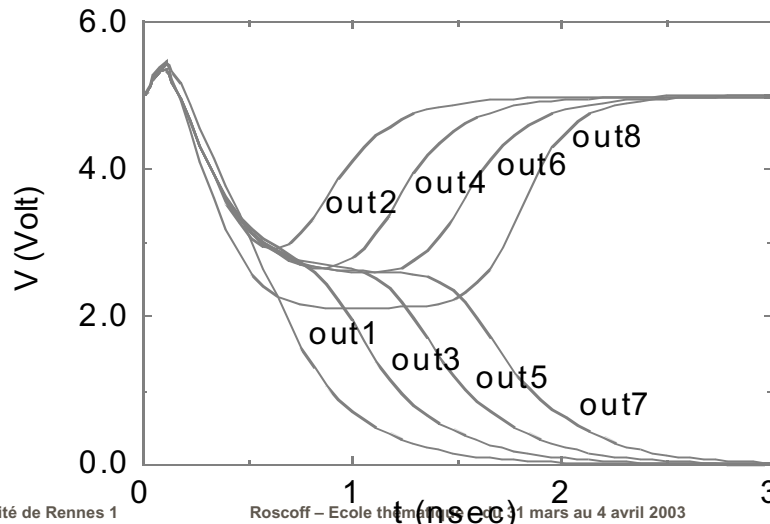
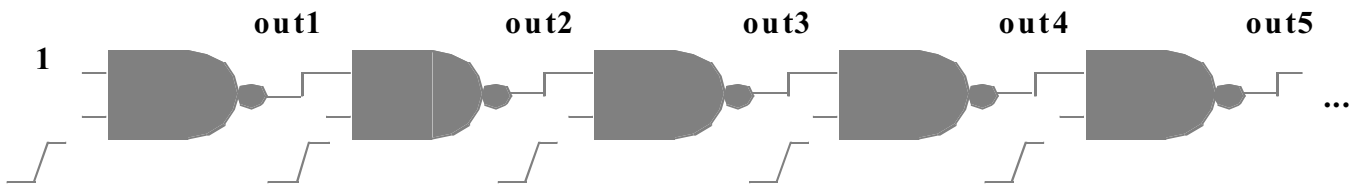
$$= 15/64$$



# Puissance dynamique

## ⌘ Influence de la propagation (glitch) :

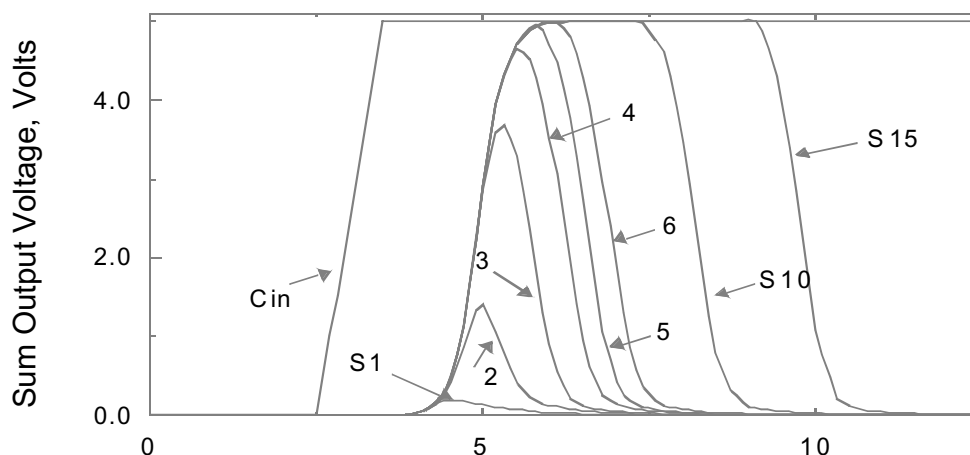
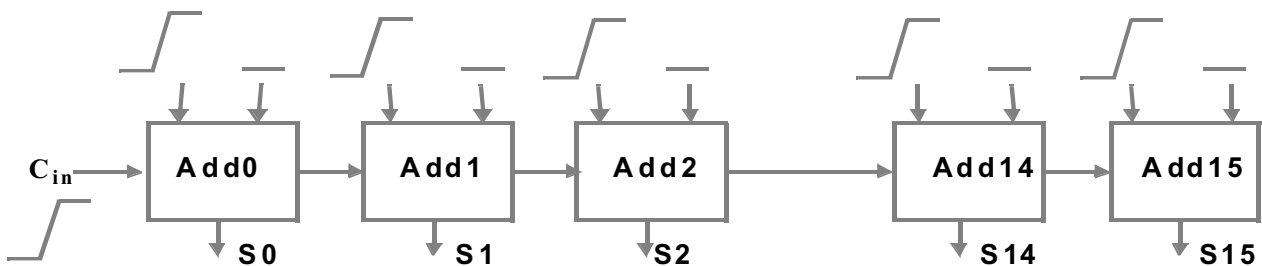
☒ si les sorties sont à 1 et que les entrées basculent de 0 à 1



# Puissance dynamique

## ⌘ Influence de la propagation (suite) :

☒ les entrées de l'additionneur passent à 1 ainsi que la retenue



## ⌘ Puissance de charge : $P_c$

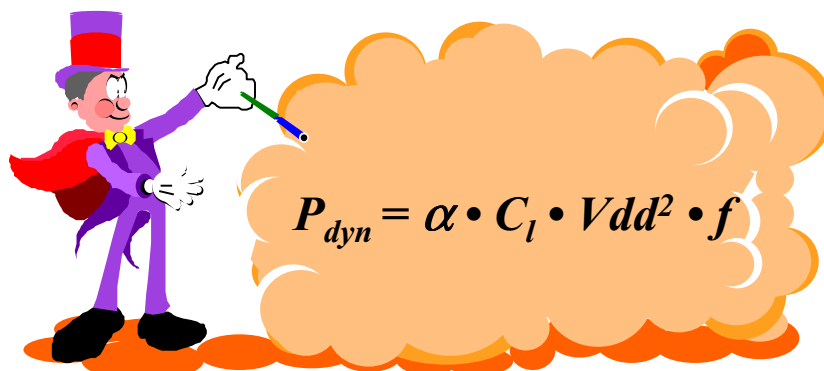
- ☑ Charge et décharge des capacités du circuit

## ⌘ Courants de court-circuit : $P_{cc}$

- ☑ Chemin de court circuit pendant la commutation des structures logiques statiques

## ⌘ Courants de fuite : $P_f$

- ☑ Jonctions, fonctionnement sous le seuil



# Degré de liberté pour diminuer la conso

## ⌘ La puissance dépend de 4 paramètres :

- ☑ fréquence
- ☑ activité
- ☑ capacité effective
- ☑ tension d'alimentation (influence quadratique !!!)

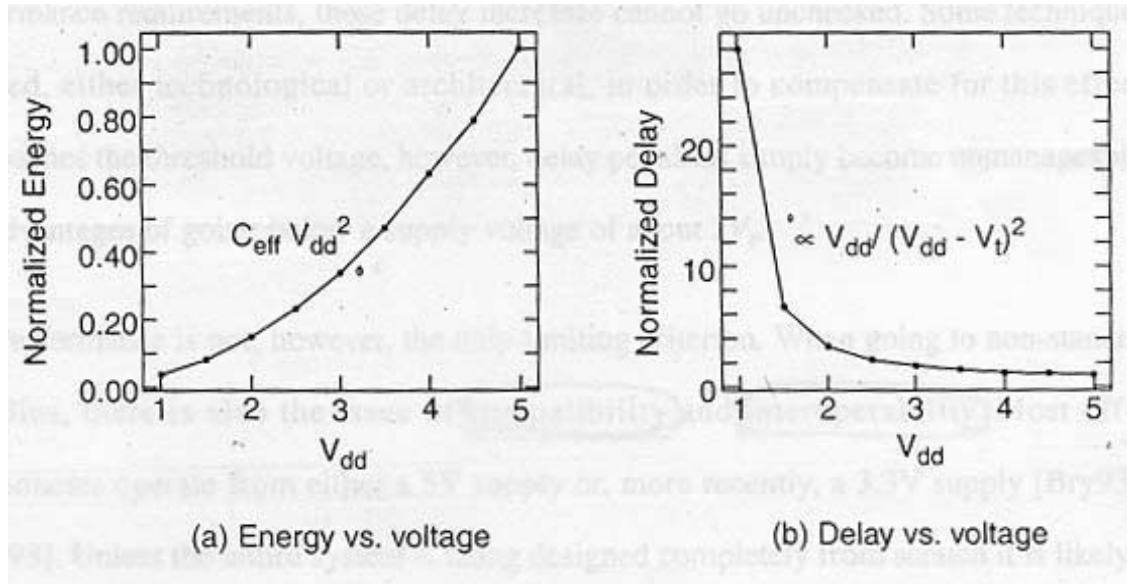
## ⌘ On peut jouer sur les 4 paramètres :

- ☑ attention, ils ne sont pas orthogonaux
- ☑ attention, ils ne peuvent être optimisés indépendamment les uns des autres

# Dépendances entre les paramètres

⌘ Exemple :

☒ tension versus délai



# Dépendances entre les paramètres

⌘ Exemple :

☒ capacité effective versus activité

$$\alpha C = C_{eff}$$

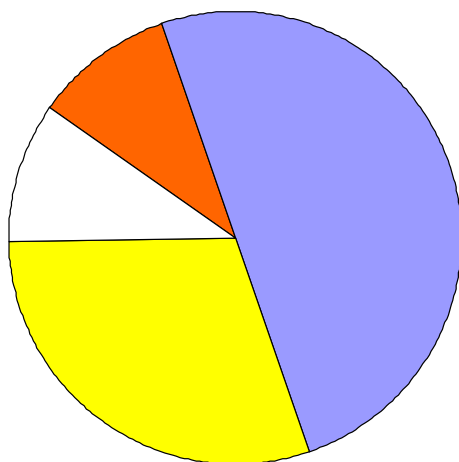


# Qu'est ce qui consomme ?

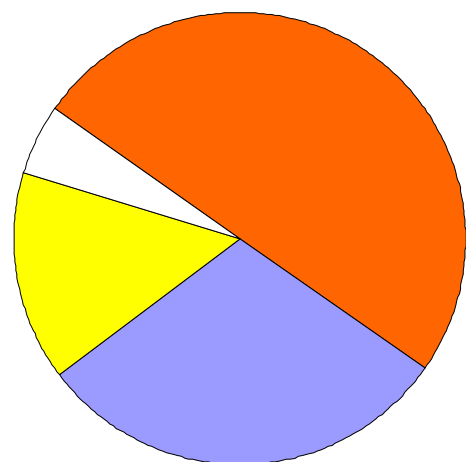
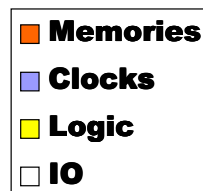
- ⌘ La logique
- ⌘ Les opérateurs
- ⌘ La mémoire
- ⌘ Les interconnexions
- ⌘ L'horloge
- ⌘ L'écran
- ⌘ Les disques
- ⌘ ...
- ⌘ ...



# Qu'est ce qui consomme ?



DEC Alpha uP

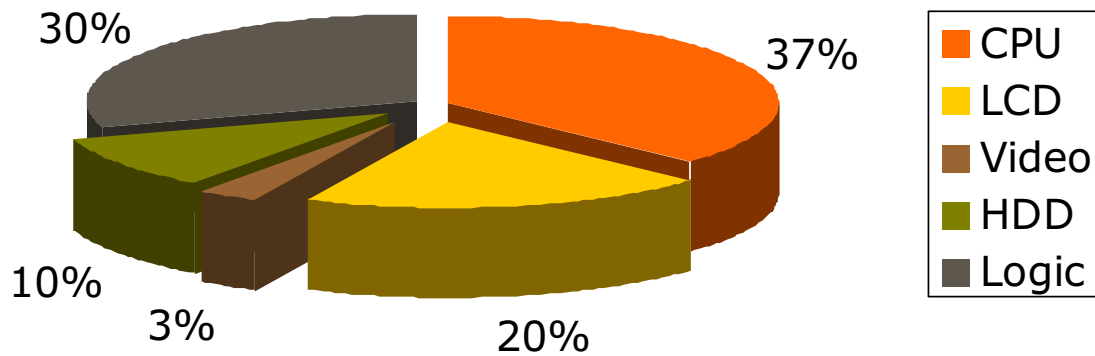


Telecom chip

# Qu'est ce qui consomme ?

## ⌘ PC Portable

☒ Puissance totale (application Word) : 19.1 W



[Source Hitachi]

# Qu'est ce qui consomme ?

## ⌘ Le CPU :

☒ les unités fonctionnelles évidemment

☒ le contrôleur :

☒ contrôle de l'exécution des instructions :

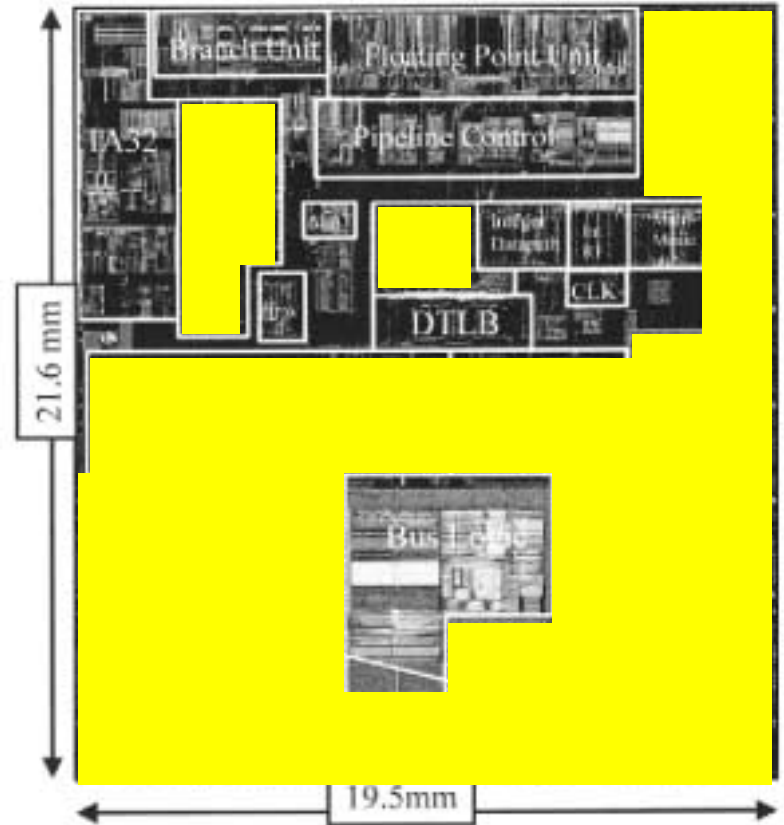
- lecture, décodage, exécution, rangement de résultat

☒ les caches qui sont de plus en plus intégrés dans le chip

# Qu'est ce qui consomme ?

## ⌘ La mémoire

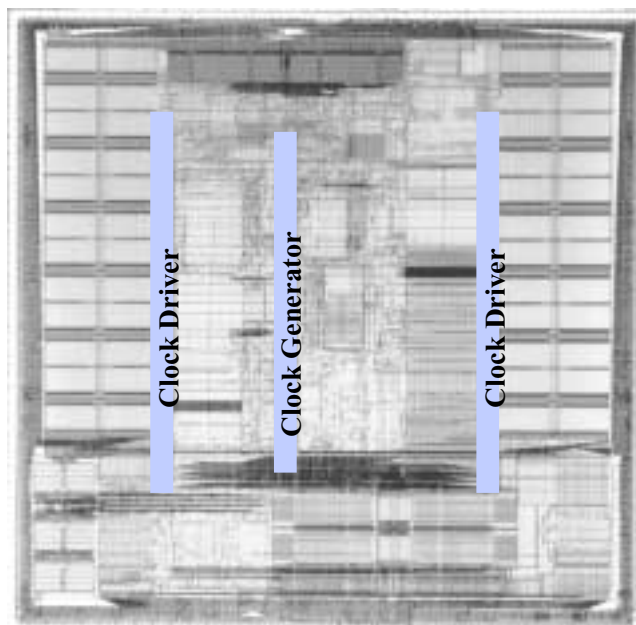
- ☒ interne (cache, scratch memory)
- ☒ externe (passage par les entrées sorties du chip)
  
- ☒ exemple de l'itanium 2
  - ☒ environ 60 à 70 % de la surface consacrée à la mémoire cache



# Qu'est ce qui consomme ?

## ⌘ L'horloge :

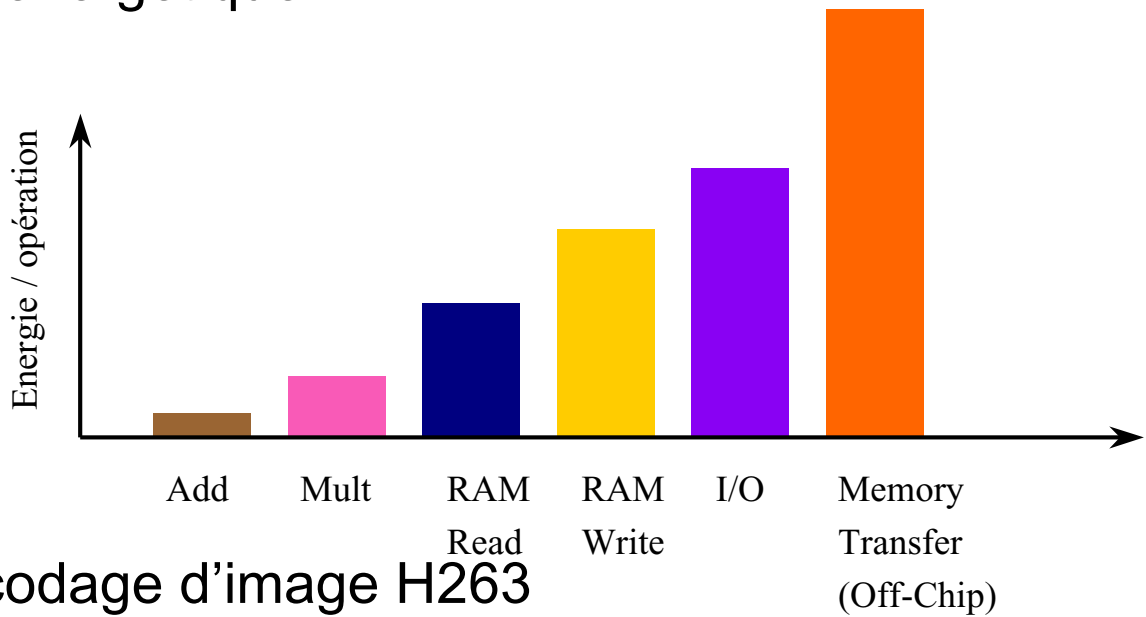
- ☒ 40-50% de la puissance dissipée par l'arbre d'horloge



**DIGITAL Corp.**  
**Processeur alpha 21164**  
**1995**  
**2.5M Portes**  
**9.3M Transistors**  
**298 mm<sup>2</sup>**  
**300 MHz**  
**64/128 Bits**  
**0.5μ**  
**60W @ 3.3V**

# Qu'est ce qui consomme ?

## ⌘ Bilan énergétique



## ⌘ E.g. codage d'image H263

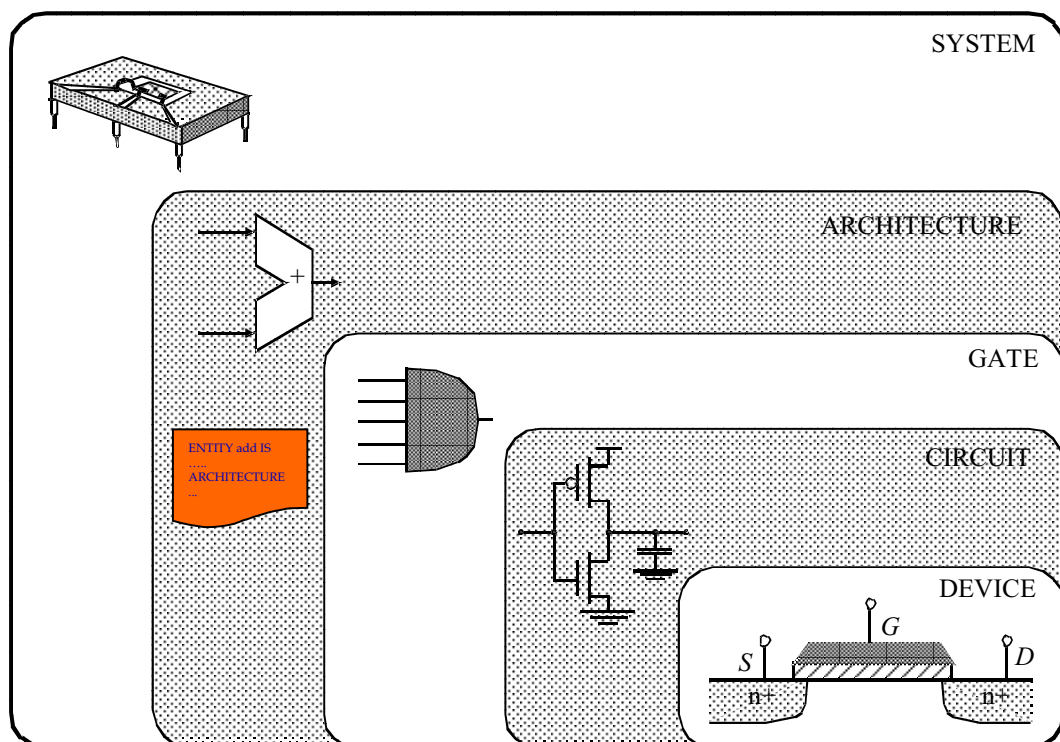
1500K Opérations et 500K Transferts mémoire

Energie transfert = 33 . Energie Opération

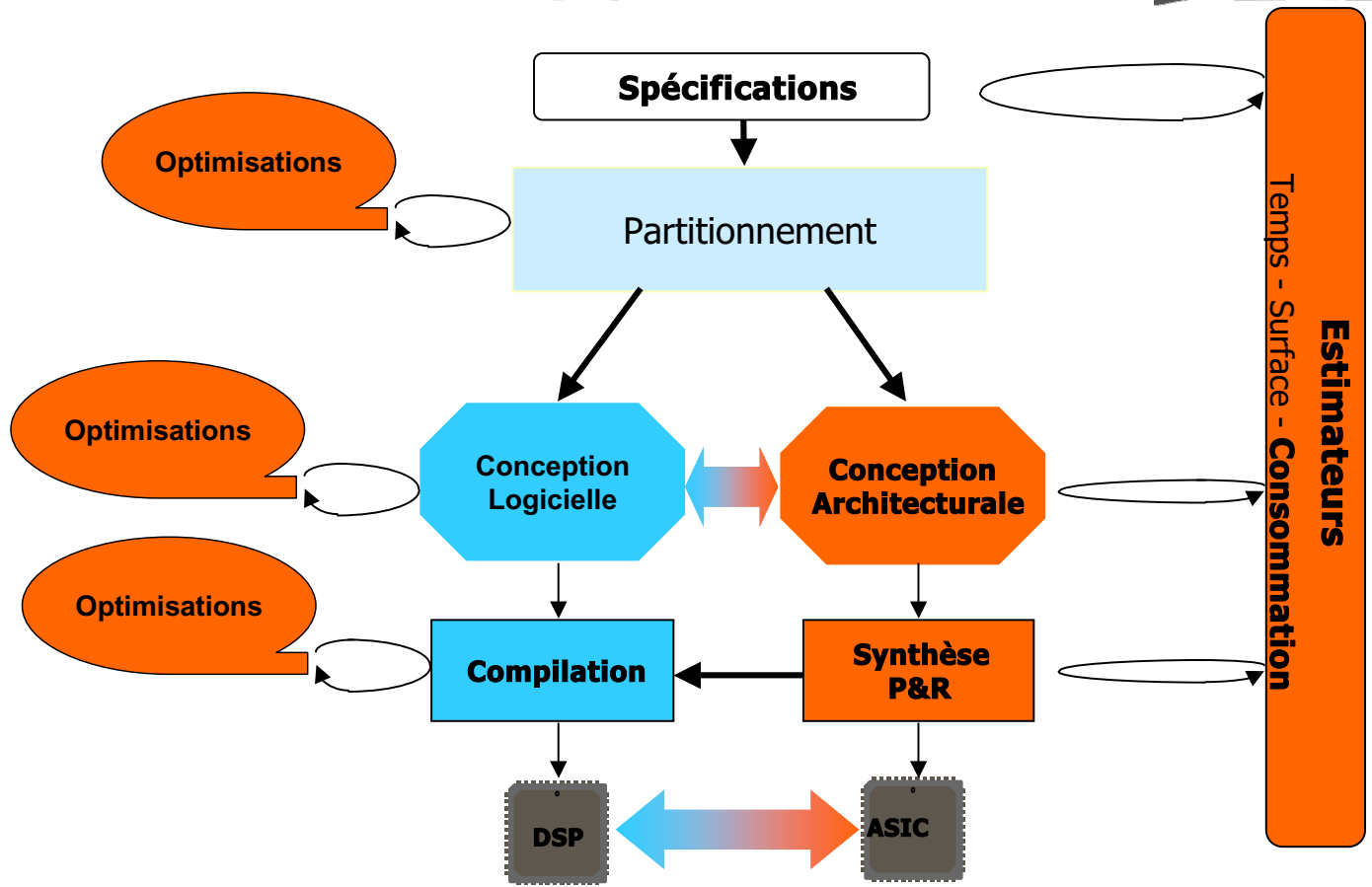
-> Puissance mémoire # 10 . Puissance opérations

# Où intervenir ?

## ⌘ Interventions possibles à tous les niveaux

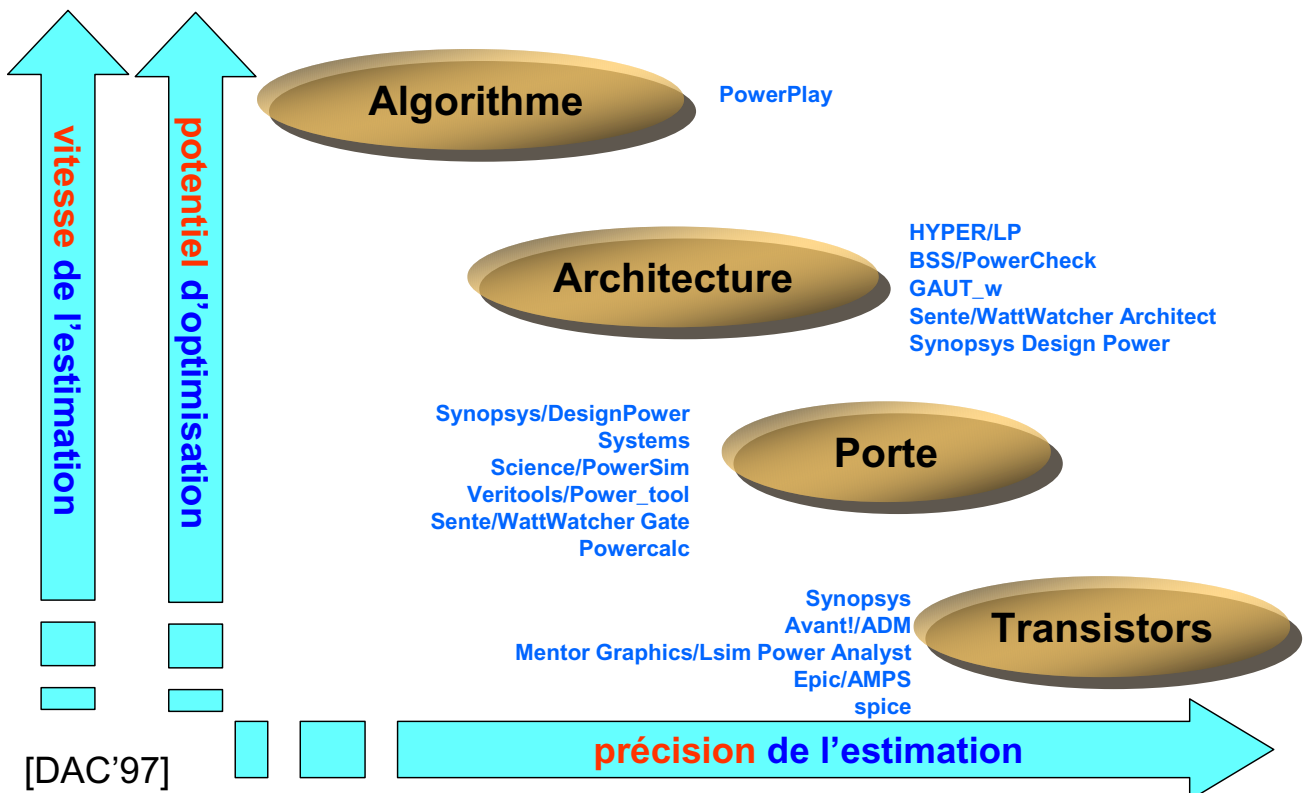


# Optimisations & estimations



# Optimisations & estimations

## ⌘ Besoins en estimateurs

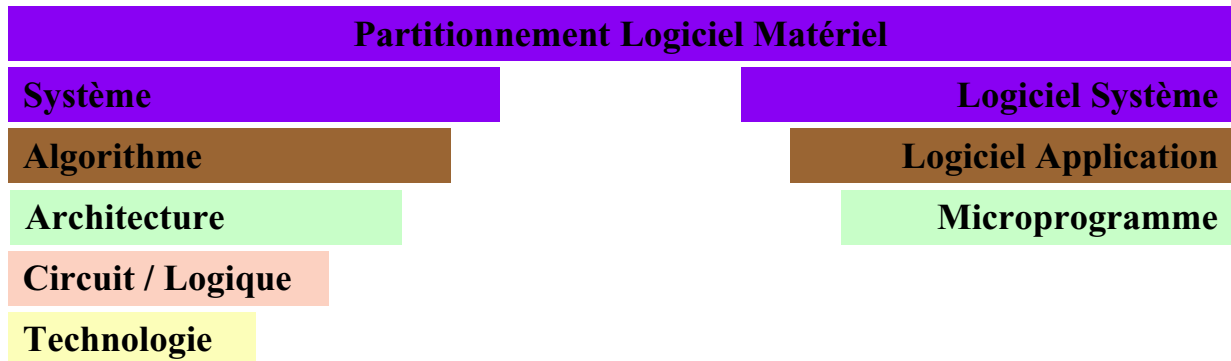


# Optimisations & estimations



⌘ Où gagner en consommation ?

Partitionnement logiciel/matériel 10-95%



Système et Algorithme 30-95%

Logiciel 10-80% (sur le système)

Architecture 10-90%

Logique 15%

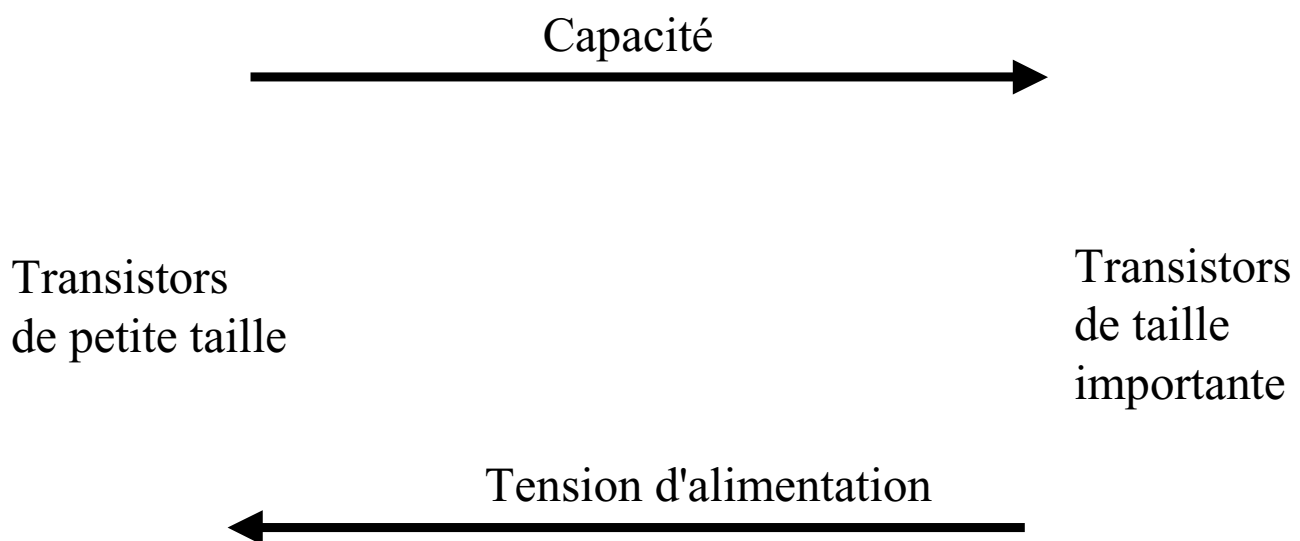
Physique 20%

Technologie 30%

# Optimisations : au niveau technologique



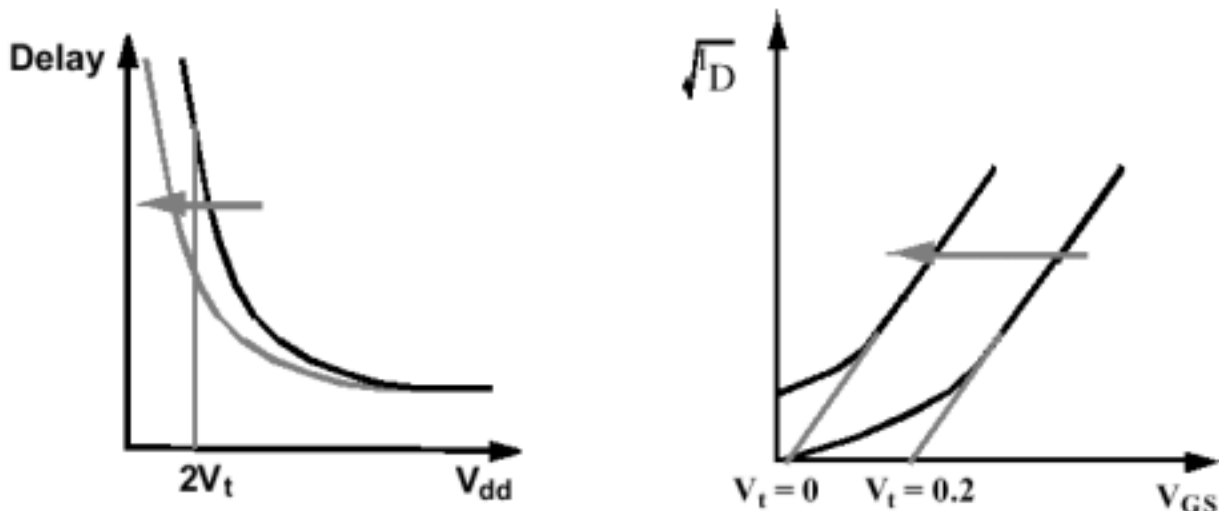
⌘ Dimensionnement des transistors :



## ⌘ Augmentation de la tension de seuil des transistors :

☒ diminution du courant de fuite :

☒ impact sur la puissance statique et dynamique



## ⌘ Abaisser la tension d'alimentation :

☒ intervient de façon quadratique

☒ de tout le circuit lorsque celui-ci est inactif ou en sous utilisation

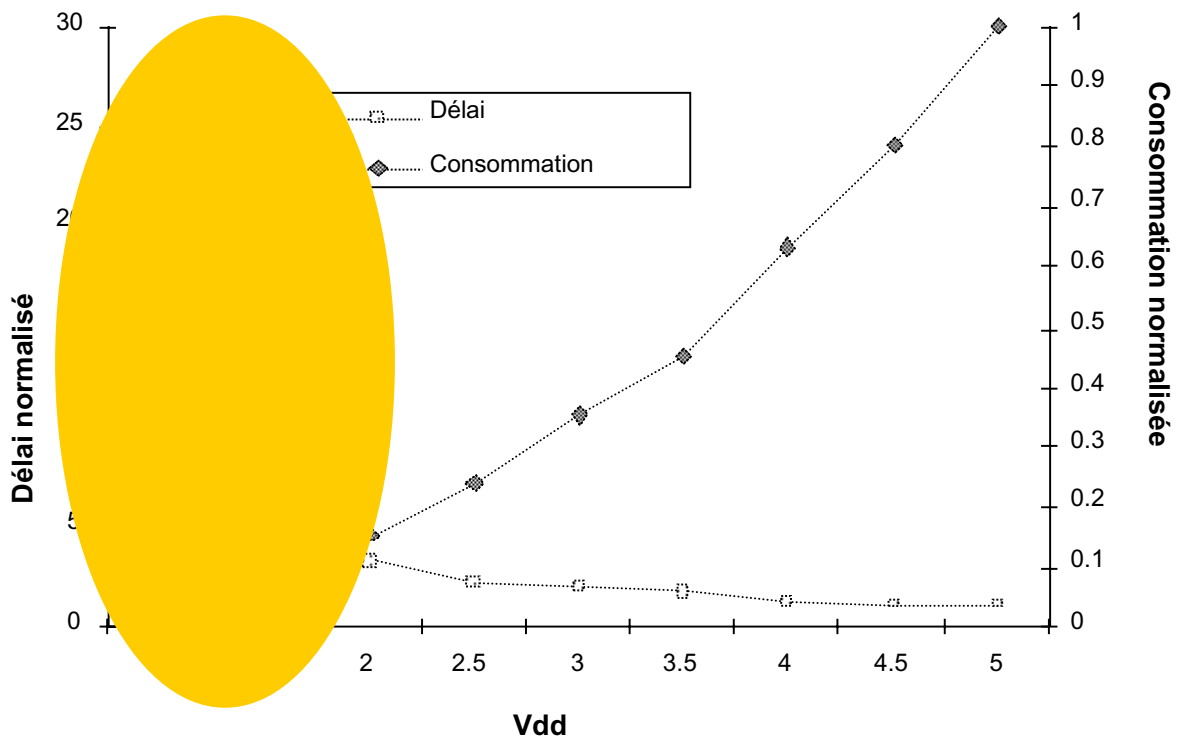
☒ d'une partie du circuit lorsque celle-ci est inutilisée ou non critique

☒ disponibilité de plusieurs tensions d'alimentation

☒ gestion dynamique de la tension

☒ attention aux interfaces avec le monde extérieur

## ⌘ Abaisser la tension d'alimentation : suite

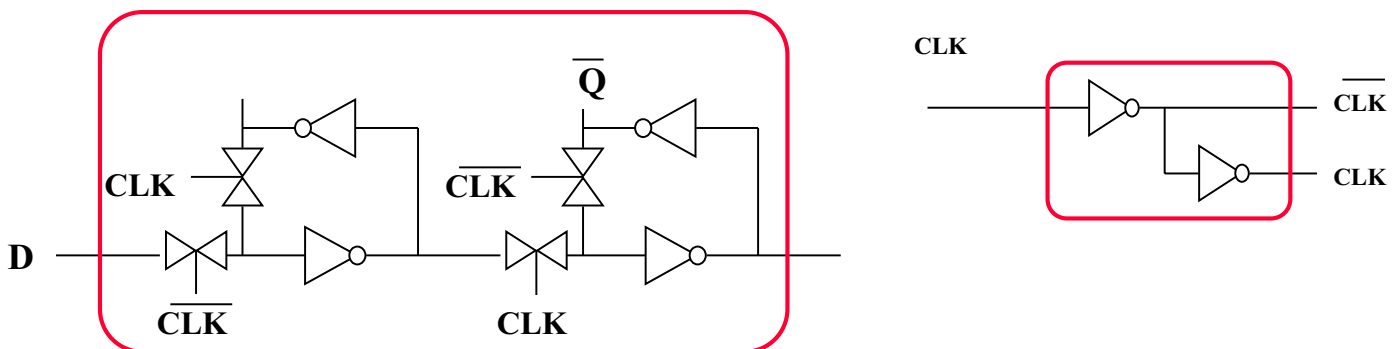


Proche de la tension de seuil

# Optimisations : au niveau logique

## ⌘ Gated clock :

- ⊞ la puissance d'une bascule dépend :
  - ⊞ de l'activité de la bascule : transitions de 0 à 1
  - ⊞ mais même sans activité, la bascule consomme



⊞ puissance :  $P_{\text{basc}} = \alpha P_{\text{dyn}} + P_{\text{horl}}$

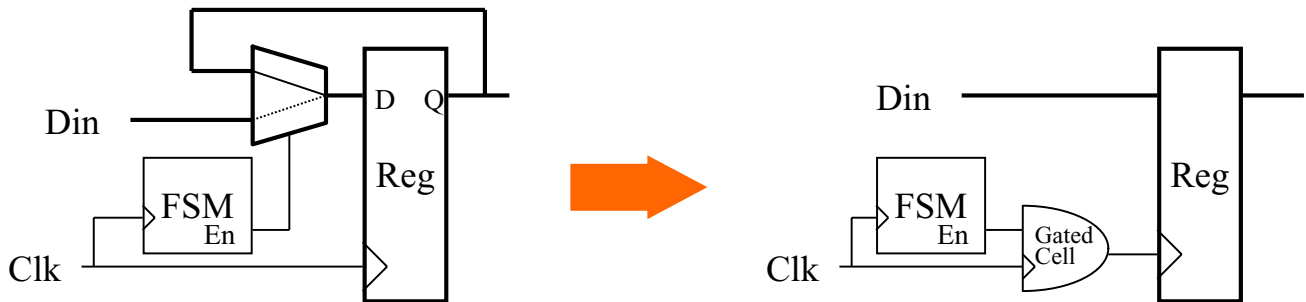
⊞ objectif du gated clock :  $P_{\text{basc}} = \alpha(P_{\text{dyn}} + P_{\text{horl}})$



# Optimisations : au niveau logique

## ⌘ Gated clock suite :

- ☑ Supprimer les commutations des blocs qui n'ont pas besoin d'être rechargés

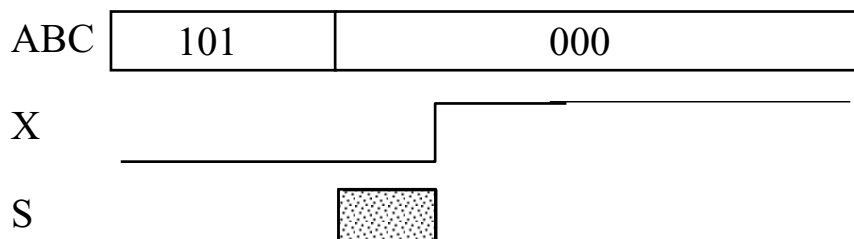
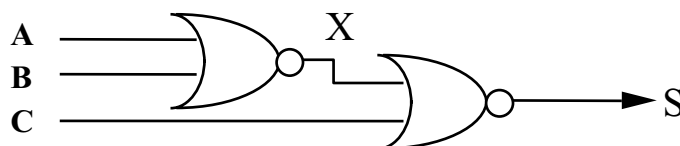


- ☑ Modification de la machine d'états
- ☑ Gain important
- ☑ Mais... Circuit non totalement synchrone

# Optimisations : au niveau logique

## ⌘ Glitch (hasards dynamiques)

- ☑ Transitions dynamiques parasites
- ☑ Consommation majeure (ET INUTILE) de courant



# Optimisations : au niveau logique

## ⌘ Elimination des glitches :

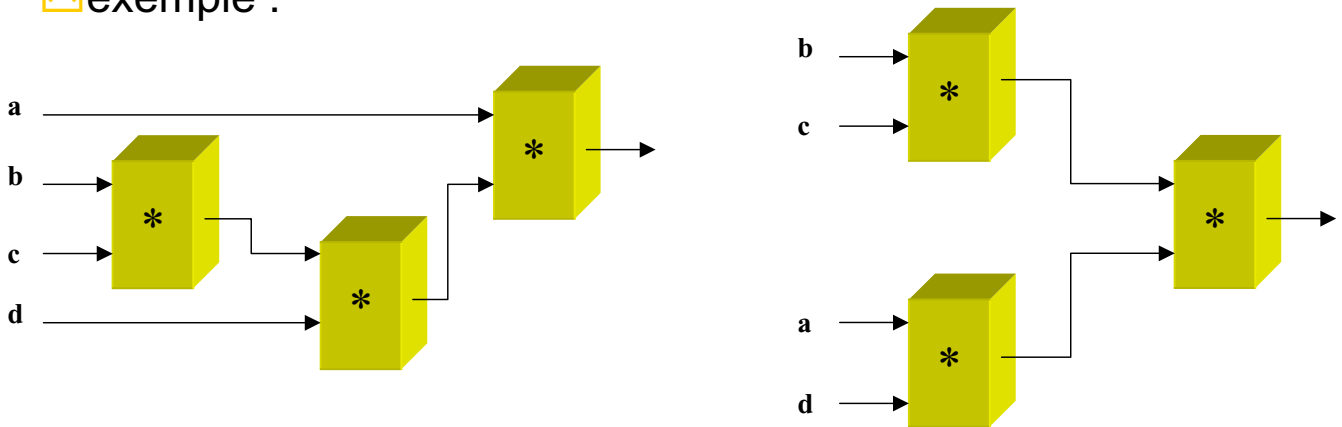
☑ par réorganisation des fonctions combinatoires lorsque cela est possible :

☑ limitation de la longueur du chemin :

- stabilité de la sortie obtenue plus rapidement
- moins de commutations

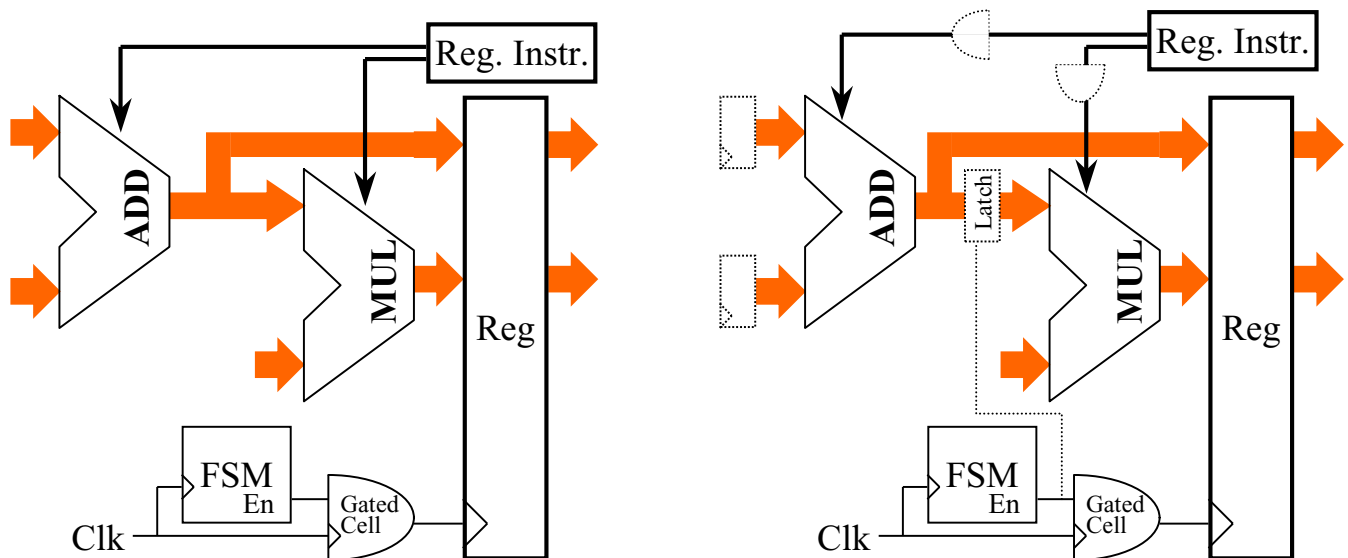
☑ limitation de la capacité effective

☑ exemple :



# Optimisations : au niveau logique

## ⌘ Activer les UF uniquement lorsque cela est nécessaire



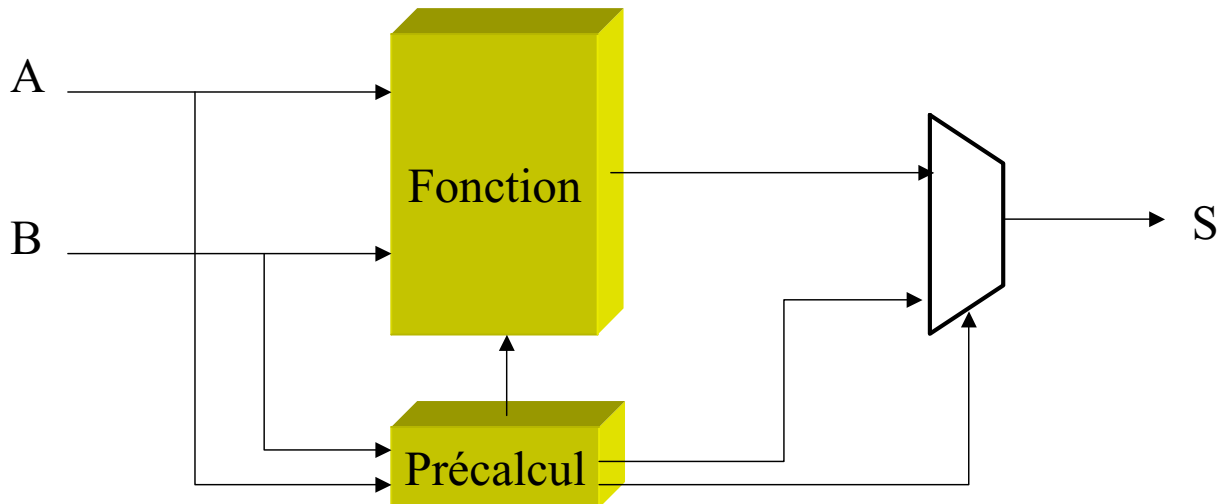
Le multiplieur consomme même si ses sorties sont inutilisées

Les entrées du multiplieur sont bloquées si ses sorties sont inutilisées  
Les instructions sont bloquées

# Optimisations : au niveau logique

## ⌘ Pré calcul :

- ☒ principe général : éviter le passage par une fonction coûteuse lorsque le résultat peut être pré calculé

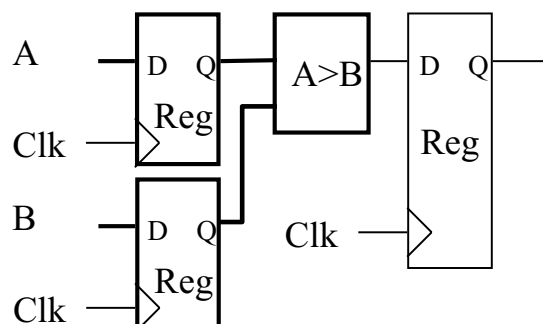


# Optimisations : au niveau logique

## ⌘ Pré calcul : exemple

- ☒ comparateur  $A > B$
- ☒ si les 2 MSB sont différents alors le résultat peut être déterminé sans soustraire A et B :
  - ☒ Si  $A[\text{MSB}] \neq B[\text{MSB}] \ \&\& \ A[\text{MSB}] == 0$  alors  $A > B$  est vraie (1)
  - ☒ Si  $A[\text{MSB}] \neq B[\text{MSB}] \ \&\& \ A[\text{MSB}] == 1$  alors  $A > B$  est faux (0)

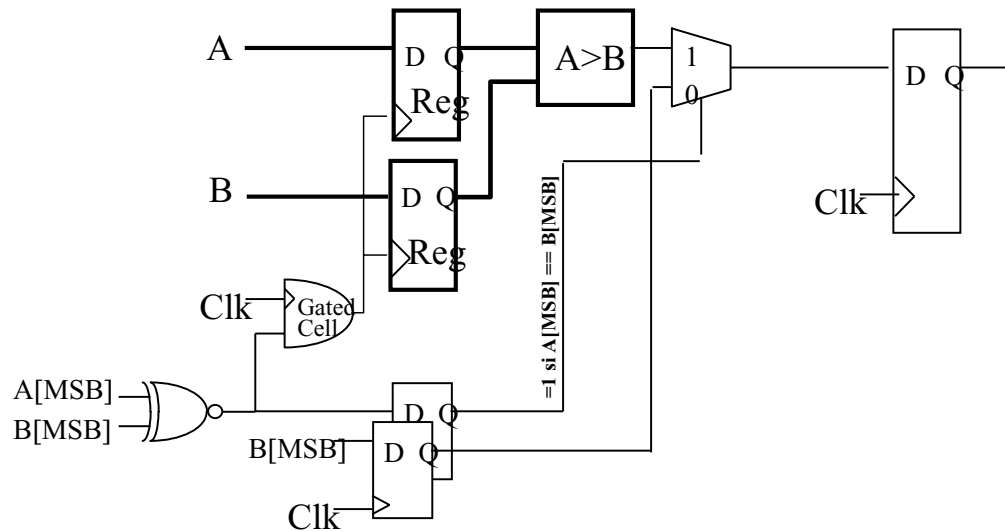
Schéma classique



# Optimisations : au niveau logique

## ⌘ Pré calcul : exemple

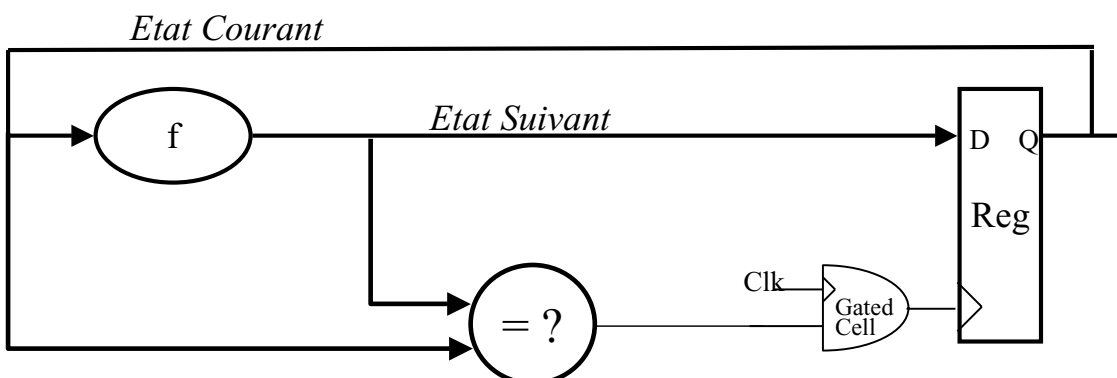
- ☒ si les 2 MSB sont égaux alors il faut faire la soustraction
- ☒ sinon, il suffit de prendre la valeur du MSB de B :
  - ☒ si MSB B == 1 alors  $A > B$  est vraie
  - ☒ sinon  $A > B$  est faux



# Optimisations : au niveau logique

## ⌘ Post calcul :

- ☒ principe général : bloquer un chargement si l'état suivant est le même que l'état courant



# Optimisations : au niveau logique

## ⌘ Post calcul :

☑ exemple

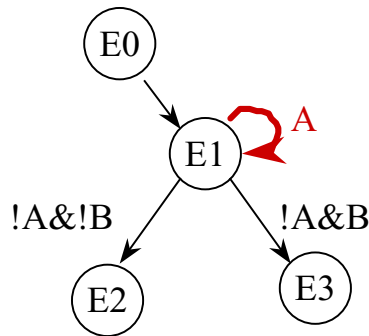
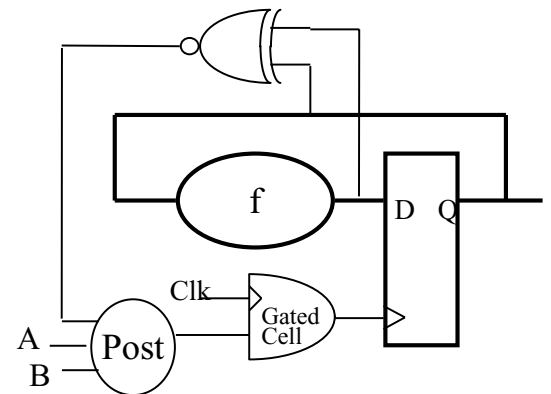
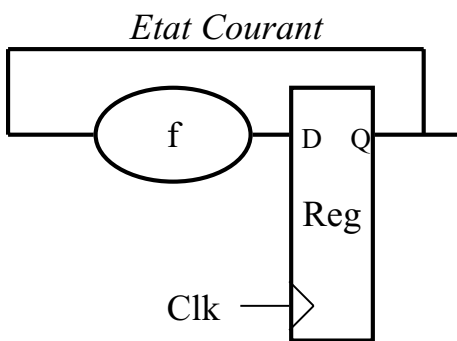


Schéma bloquant le chargement du même état

Schéma classique



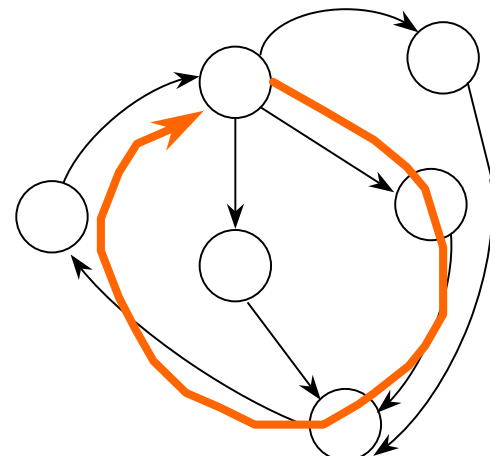
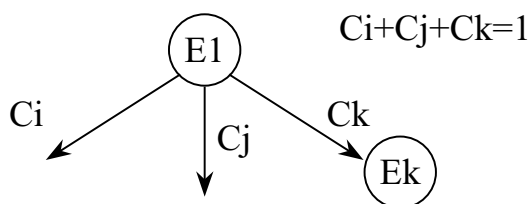
# Optimisations : au niveau logique

## ⌘ Les machines d'états :

☑ Code pour les états d'une machines

**Gray**  
- Activité

**Binaire**  
+ Activité



### Coder le graphe avec le minimum de transitions

Si Prob(Ck) importante Alors  
les codes de E1 et Ek doivent différer  
avec un minimum de positions  
(distance de Hamming)

**Chemin le plus probable**  
A coder avec un Min de transitions

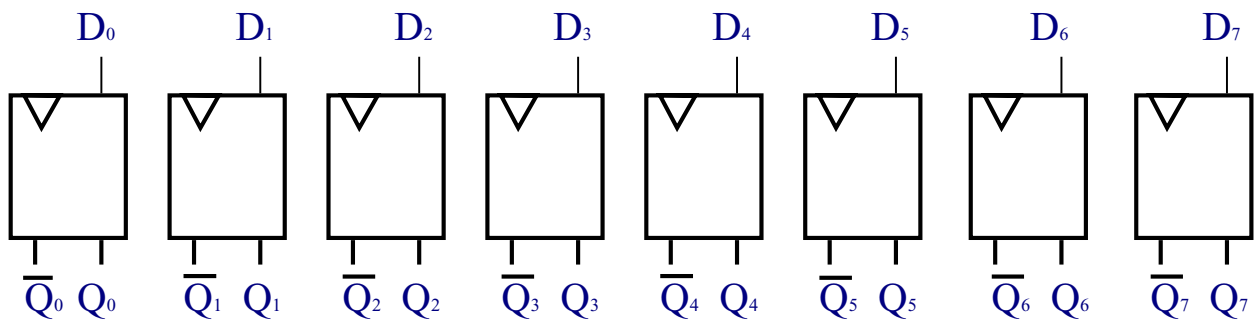
# Optimisations : Dépendance aux données



⌘ e.g. Registre 8 bits

⌘ Consommation d'une bascule :

⊠ 9.5µW/MHz par transition de 0 à 1 à la sortie



# Optimisations : Dépendance aux données

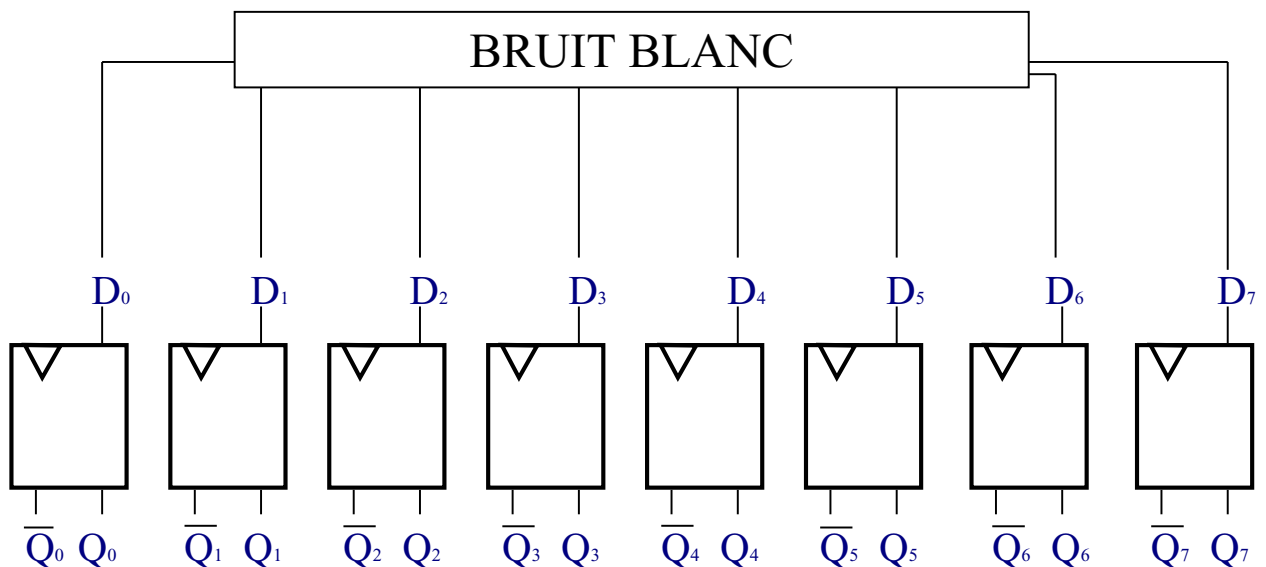


⌘ Signal aléatoire :

⊠ Environ 4 bascules commutent pour un signal aléatoire en entrée :

⊠ deux bascules commutent de 1 à 0 et deux de 0 à 1 :

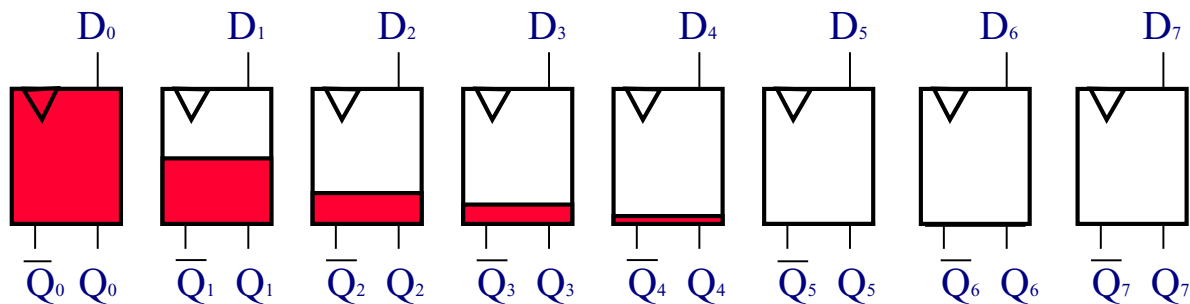
- donc  $2 * 9,5 = 19 \mu\text{W}/\text{MHz}$



# Optimisations : Dépendance aux données

## ⌘ Comptage binaire en complément à 2

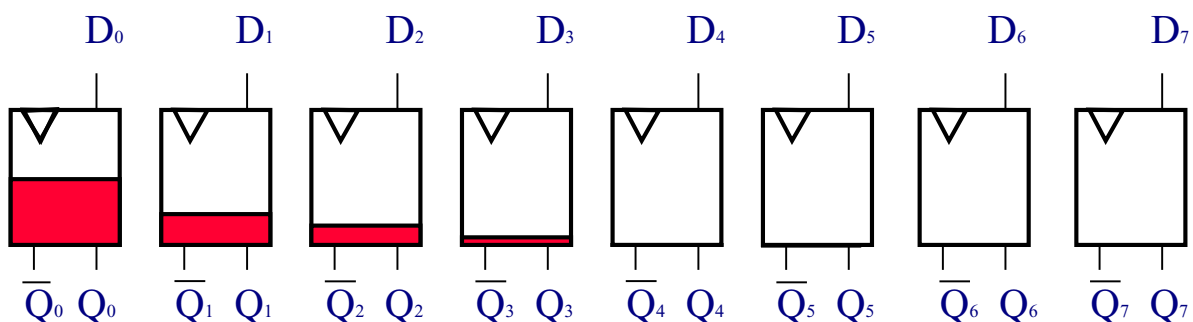
- ☒ Proba transition =  $PT(b_0) + PT(b_1) + PT(b_2) + \dots$  0 0 0
- ☒  $1 + 1/2 + 1/4 + 1/8 \dots = 2$  0 0 1
- ☒ conso =  $2 * 9,5 / 2 = 9,5 \mu\text{W}/\text{MHz}$  0 1 0
- 0 1 1
- 1 0 0
- 1 0 1
- 1 1 0
- 1 1 1



# Optimisations : Dépendance aux données

## ⌘ Comptage binaire en codage gray

- ☒ Proba transition =  $PT(b_0) + PT(b_1) + PT(b_2) + \dots$  0 0 0
- ☒  $1/2 + 1/4 + 1/8 \dots = 1$  0 0 1
- ☒ conso =  $1 * 9,5 / 2 = 4,75 \mu\text{W}/\text{MHz}$  0 1 1
- 0 1 0
- 1 1 0
- 1 1 1
- 1 0 1
- 1 0 0



# Optimisations : Dépendance aux données



⌘ Dépendance aux données ???

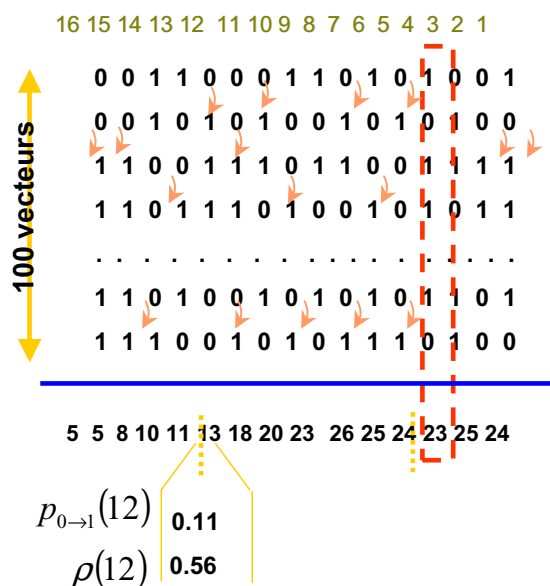
Signal	Bruit blanc	Signal de comptage binaire complément à 2	Signal de comptage binaire gray
Consommation	19 $\mu$ W/MHz	9,5 $\mu$ W/MHz	4,75 $\mu$ W/MHz

- ⌘ D'où l'importance d'avoir un modèle de signal d'entrée
- ⌘ et de pouvoir propager ce modèle dans l'architecture

# Optimisations : Dépendance aux données

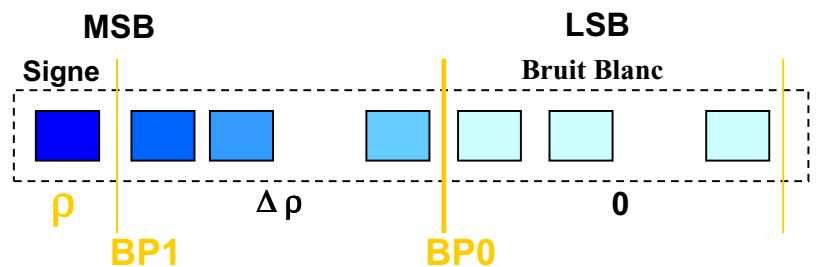


[PREECE81]



## Dual Bit Type

[LANDMAN96]



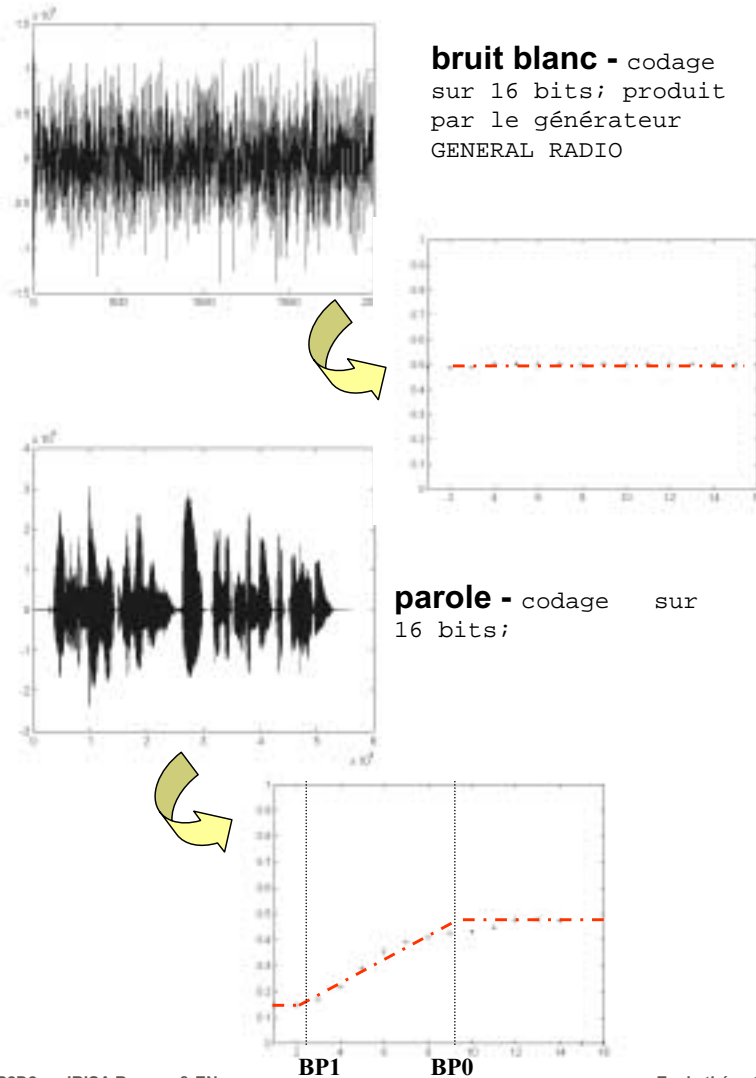
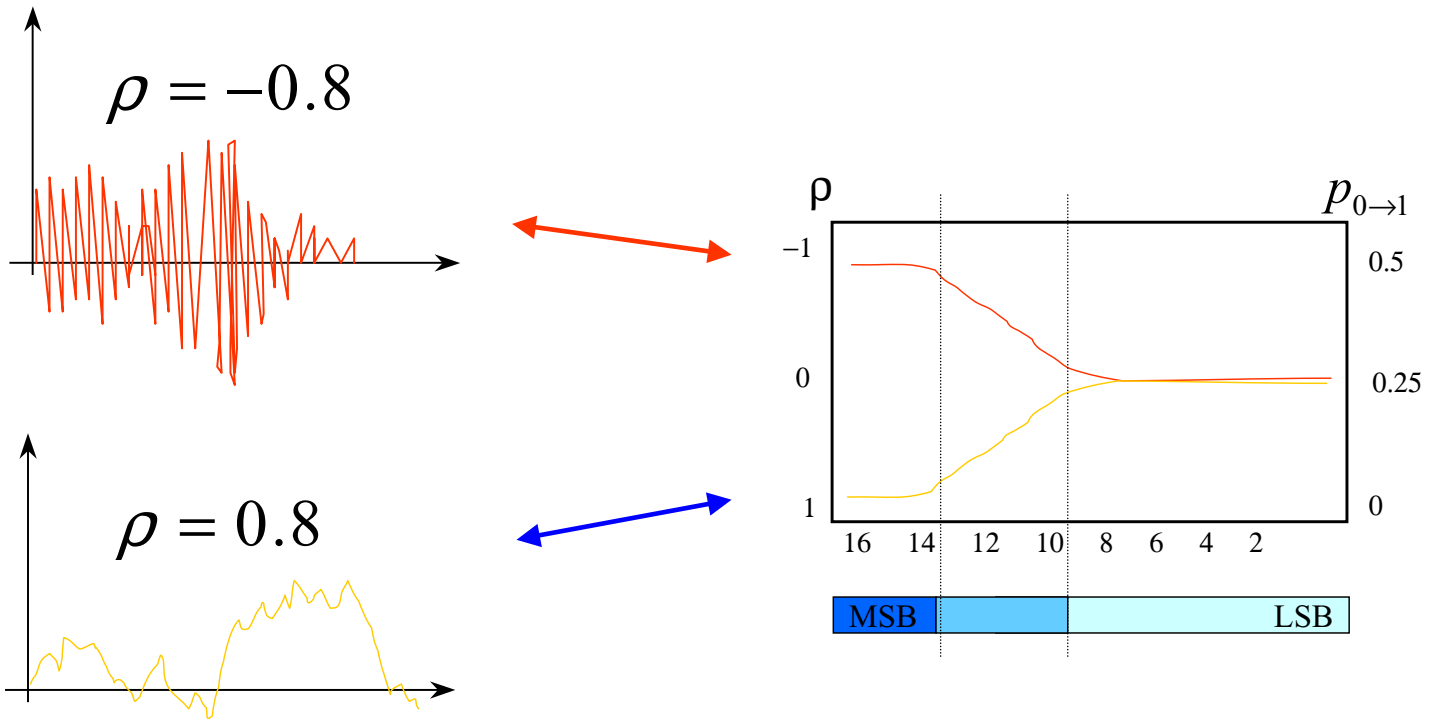
⌘ Signaux audio, vidéo, radio,

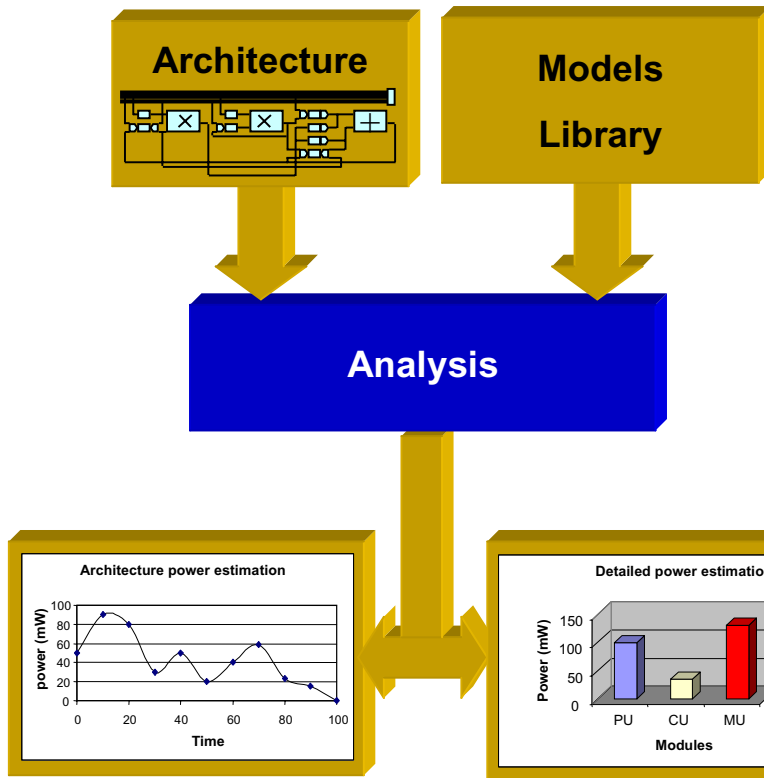
...

$$\rho = 1 - 4 \cdot p_{0 \rightarrow 1}$$



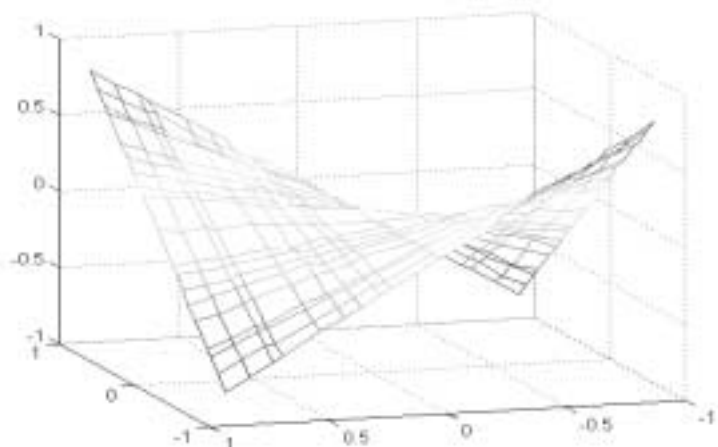
# Optimisations : Dépendance aux données





- Modèle probabiliste du signal (DBT)
- Propagation des propriétés du signal dans l'archi
- Estimation conso, : traitement, mémoire, contrôle

## ⌘ Propagation du signal dans l'architecture

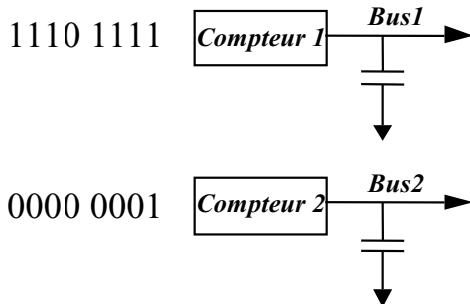


# Optimisations : au niveau architectural

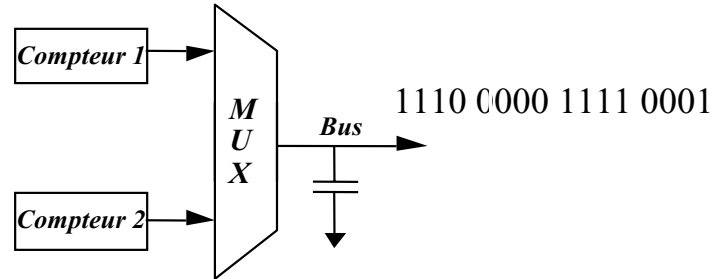


## ⌘ Le partage de ressources augmente l'activité

- ☒ Destruction des corrélations entre signaux



a) Implémentation parallèle



b) Multiplexage temporel

- ☒ Effet de la réutilisation des ressources :

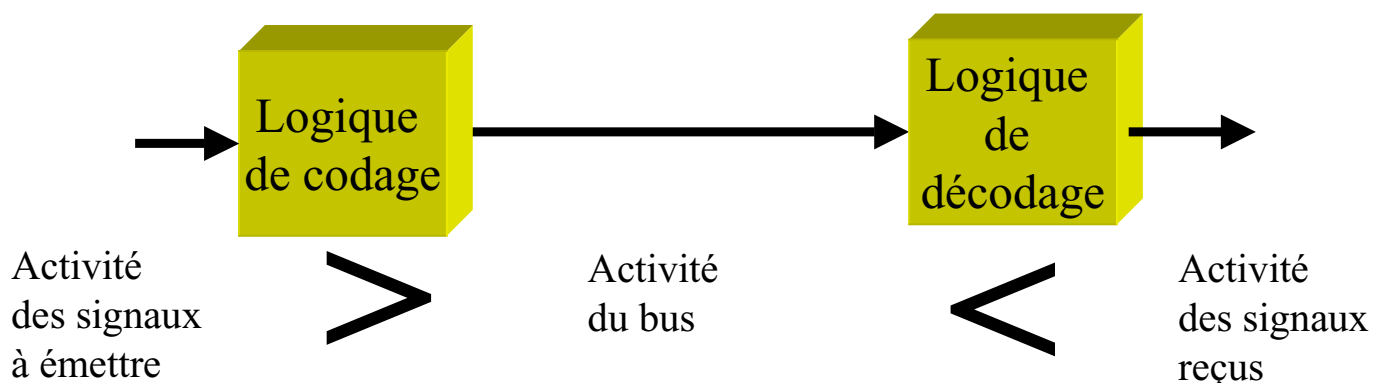
- ☒ non respect des flots de données
- ☒ augmentation de l'activité
- ☒ augmentation de la conso

# Optimisations : au niveau architectural



## ⌘ Encodage de bus pour diminuer l'activité

- ☒ e.g. liaison mémoire ou cache (instruction ou données) vers CPU
- ☒ objectif : diminuer le nombre de transitions entre les valeurs véhiculées :
  - ☒ codages :
    - activité binaire > activité gray

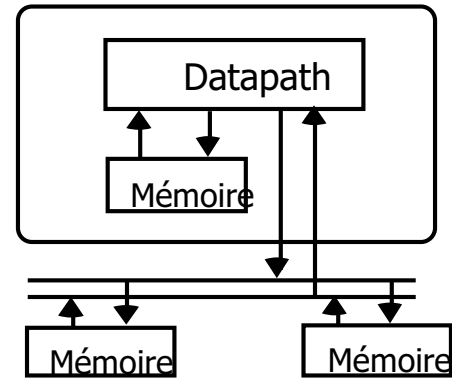


# Optimisations : au niveau architectural

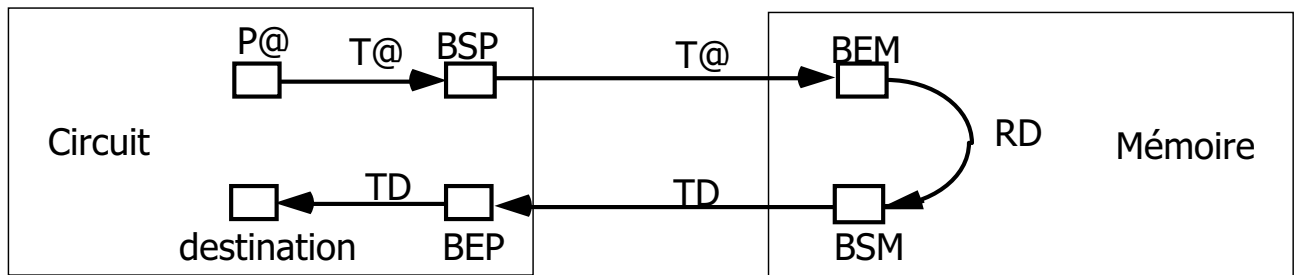


## ⌘ Les mémoires :

### ☑ Modélisation générale



### ☑ modélisation de la consommation lors d'un accès mémoire



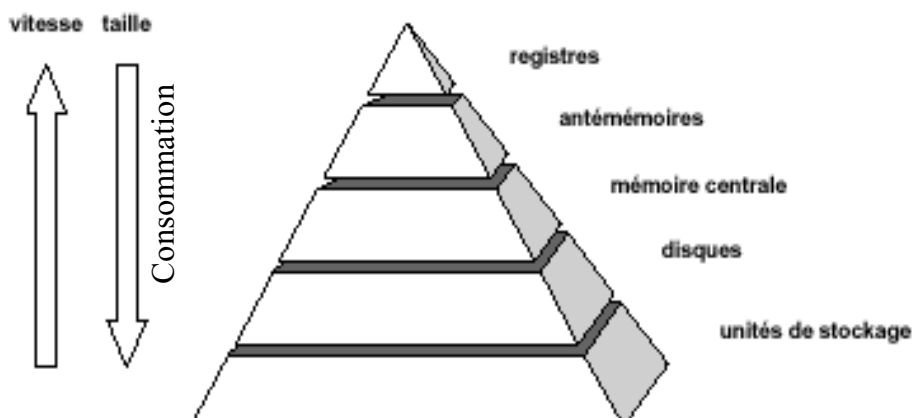
# Optimisations : au niveau architectural



## ⌘ Organisation mémoire :

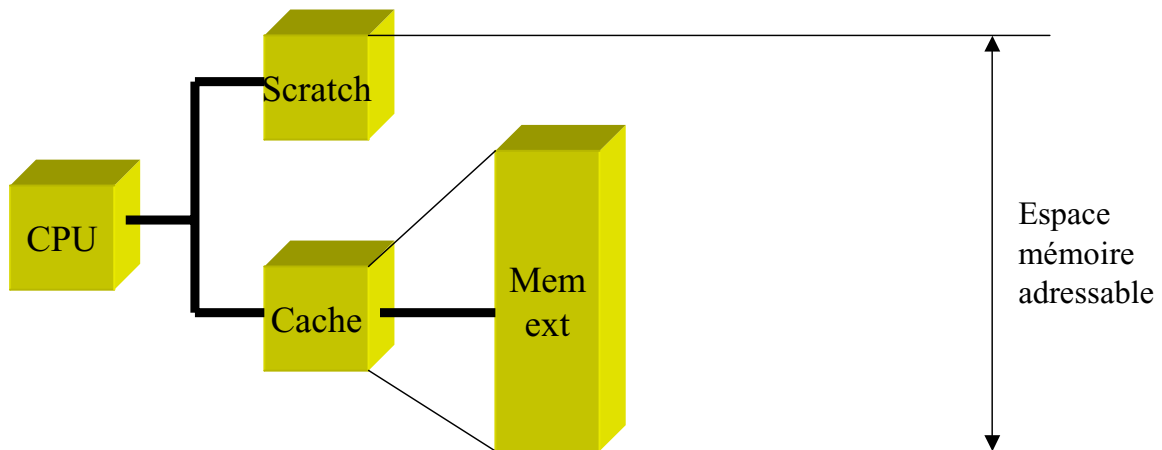
### ☑ Constats et motivations :

- ☑ croissance importante du coût engendré par la mémorisation
- ☑ la mémoire est très souvent le goulot d'étranglement du système
- ☑ modèle sans hiérarchisation :
  - mal adapté, beaucoup d'accès à une grosse mémoire
  - nécessite de disposer de mémoires importantes et avec temps d'accès court



## ⌘ Organisation mémoire :

- ☒ placement des données les plus souvent utilisées dans des mémoires internes (cache)
- ☒ minimisation de la taille mémoire totale par réutilisation des points mémoires
- ☒ mise en place de buffers ou de scratchpad pour limiter les accès aux mémoires



## ☒ techniques :

- ☒ exploiter les profils des accès aux données
  - localités spatiales et temporelles des données
- ☒ exploiter le déterminisme de l'application :
  - nombre d'accès à chaque données
  - dates des accès
  - sens des accès
  - décision de placement (dans le scratch-pad ou en mémoire principale)

# Optimisations : au niveau architectural



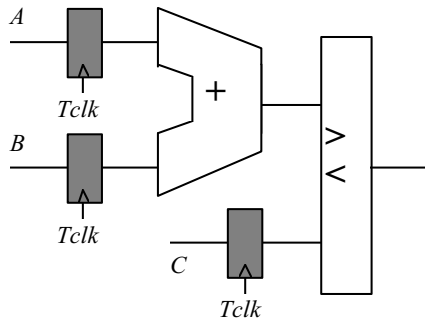
⌘ Compenser la réduction de V<sub>dd</sub> par des optimisations architecturales ou logiques

☑ Exemple : architecture 16 bits (décodage de Viterbi)

$$P_{\text{ref}} = C_{\text{ref}} \cdot V_{\text{ref}}^2 \cdot F_{\text{ref}} = 14.7 \text{ mW}$$

$C_{\text{ref}}$  déterminé par simulation (CI = 31pF)

$$F_{\text{ref}} = 1 / 40\text{ns} \quad V_{\text{ref}} = 5\text{V}$$



$$\text{Surface} = 788 \times 555 \mu\text{m}^2 = 0.44 \text{ mm}^2$$

3823 Transistors  
56 mm d'interconnexions

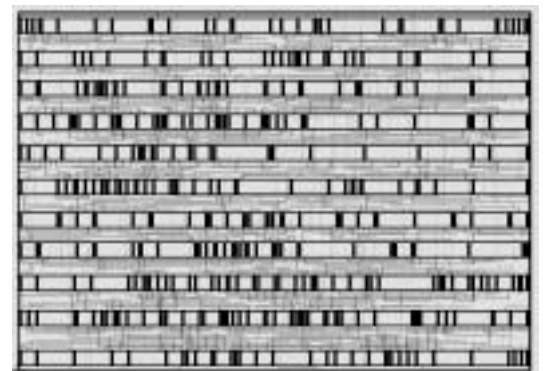
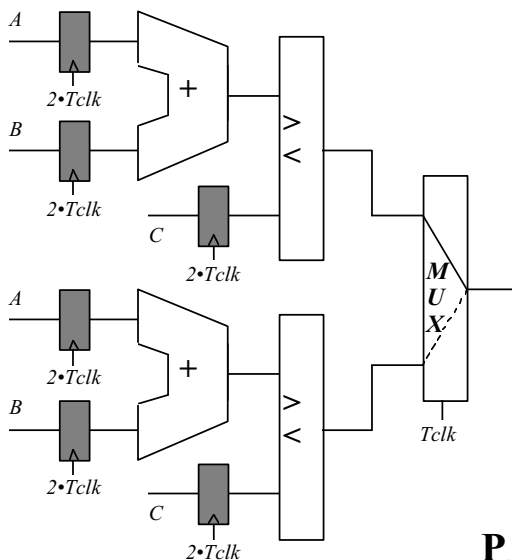
# Optimisations : au niveau architectural



⌘ Parallélisation de l'architecture

$$\text{Surface} = 1130 \times 768 \mu\text{m}^2 = 0.87 \text{ mm}^2$$

7638 Transistors  
111 mm d'interconnexions



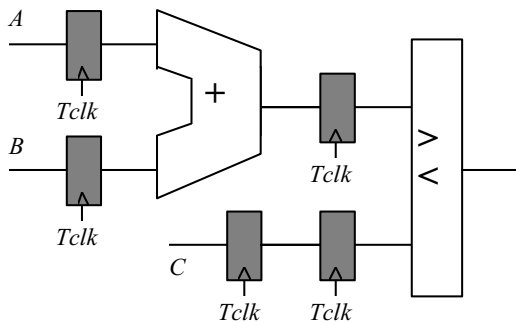
$$P_{\text{parallèle}} = (2.15 \cdot C_{\text{ref}}) \cdot (0.58 \cdot V_{\text{ref}})^2 \cdot (0.5 \cdot F_{\text{ref}})$$

$$\# 0.36 \cdot P_{\text{ref}} = 5.3 \text{ mW}$$

$$C_{\text{parallèle}} = 2.15 \cdot C_{\text{ref}}$$

$$F_{\text{parallèle}} = 1 / 80\text{ns} \quad V_{\text{dd parallèle}} = 2.9\text{V}$$

## ⌘ Pipelinage de l'architecture



$$P_{\text{pipeline}} = (1.15 \cdot C_{\text{ref}}) \cdot (0.58 \cdot V_{\text{ref}})^2 \cdot F_{\text{ref}}$$

$$\# 0.39 \cdot P_{\text{ref}} = \mathbf{5.7 \text{ mW}}$$

$$C_{\text{pipeline}} = 1.15 \cdot C_{\text{ref}}$$

$$F_{\text{pipeline}} = F_{\text{ref}}$$

$$V_{\text{dd pipeline}} = 2.9\text{V}$$

## ⌘ Parallélisation ou pipelinage ?

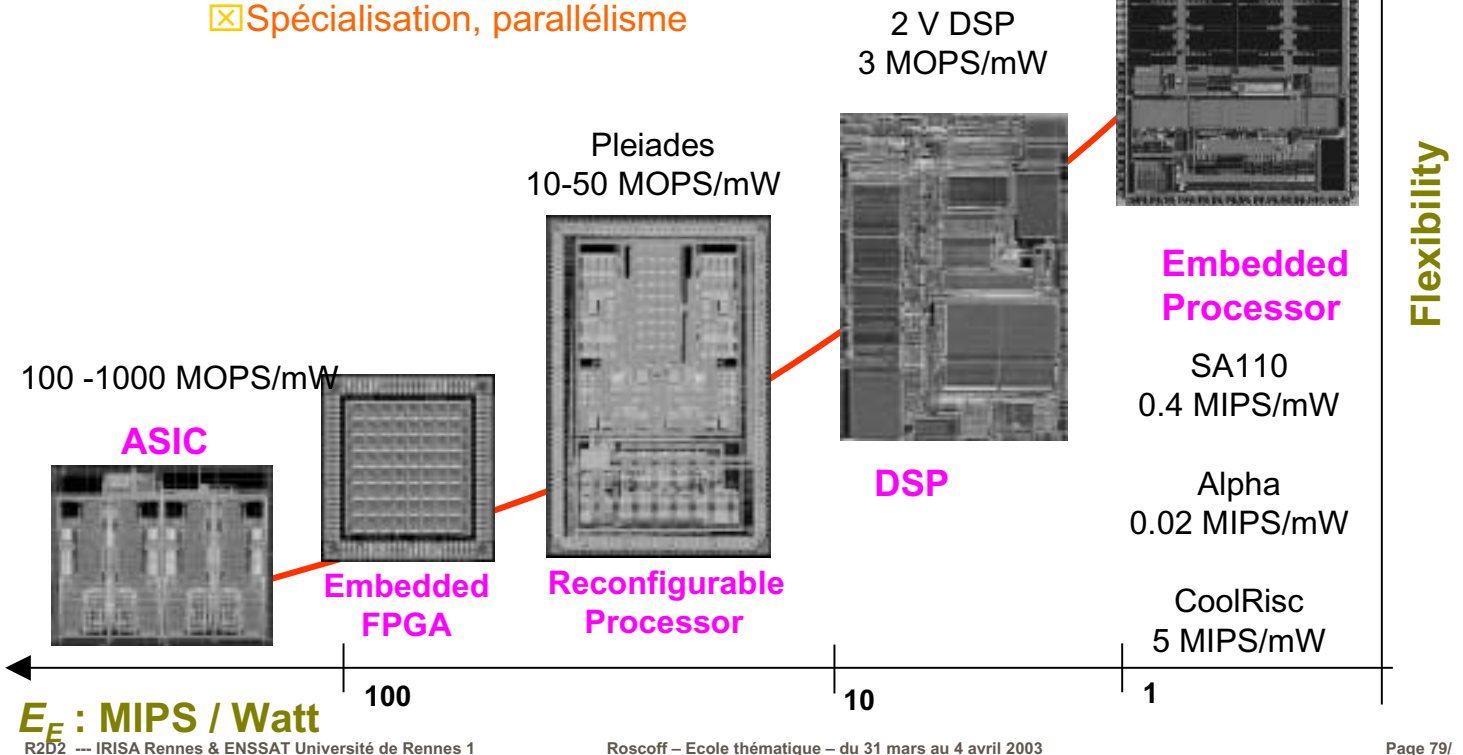
Type d'Architecture	Vdd	Surface	Puissance
Référence	5V	1	1
Parallèle	2.9V	1.98	0.36
Pipeline	2.9V	1.1	0.39
Parallèle/Pipeline	2.0V	2.2	0.2

# Alternatives architecturales

## ⌘ Compromis performance - flexibilité

⊞ Efficacité énergétique

⊞ Spécialisation, parallélisme



# Alternatives architecturales

⌘ Il ne faut plus penser uniquement aux performances

⌘ Efficacité énergétique

⊞ Nb Opérations Utiles / Energie Utilisée = NbOp / J  
= MOPS/mW

⊞ MOPS = Millions d'Opérations Par Seconde

⊞ MOPS/mW =  $(F_{clk} \cdot N_{op}) / (N_{op} \cdot A_{op} \cdot C_{sw} \cdot V_{dd}^2 \cdot F_{clk})$   
=  $1 / (A_{op} \cdot C_{sw} \cdot V_{dd}^2)$

⊞ Nop : Nb d'opérateurs

⊞ Aop : Surface moyenne d'un opérateur = Achip/Nop

⊞ Csw : Capacité efficace par mm<sup>2</sup>

⊞ dépend de la spécialisation du chip !!!

⌘ Efficacité en surface

⊞ MOPS/mm<sup>2</sup>



## ⌘ Dart : architecture reconfigurable

### ☒ objectifs :

- ☒ forte puissance de calcul
- ☒ faible consommation

### ☒ principes :

#### ☒ recherche d'un bon compromis performance / flexibilité :

- reconfiguration matérielle : exécution des cœurs de boucles
- reconfiguration logicielle : exécution d'instructions classiques

#### ☒ remplacement du contrôle d'instructions d'exécution en instructions de reconfiguration

### ☒ collaboration ENSSAT/UBO-AS/STMicroelectronics

## ⌘ Architecture :

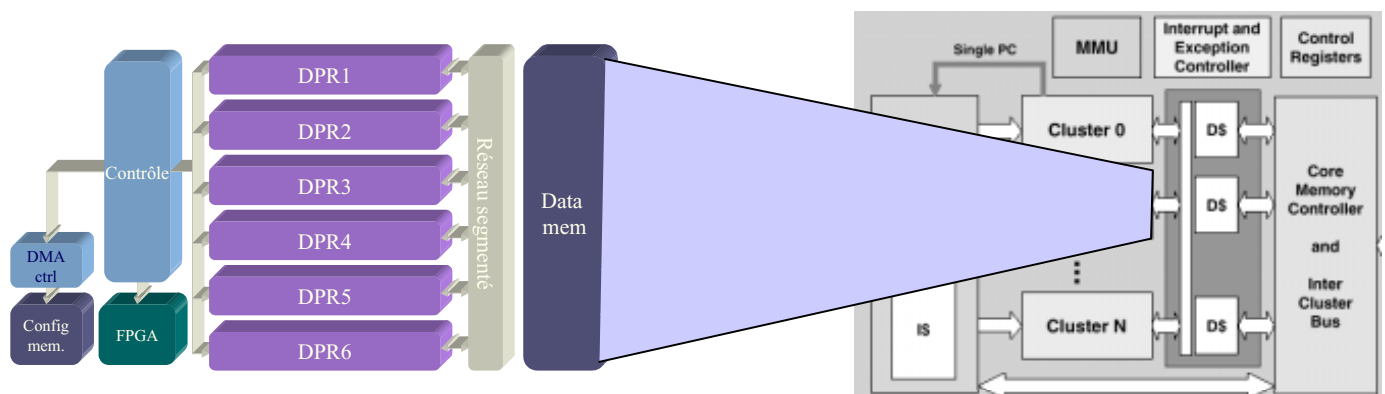
### ☒ Reconfiguration multigrain dynamique

- ☒ Fonctionnel (DPR), logique (FPGA)

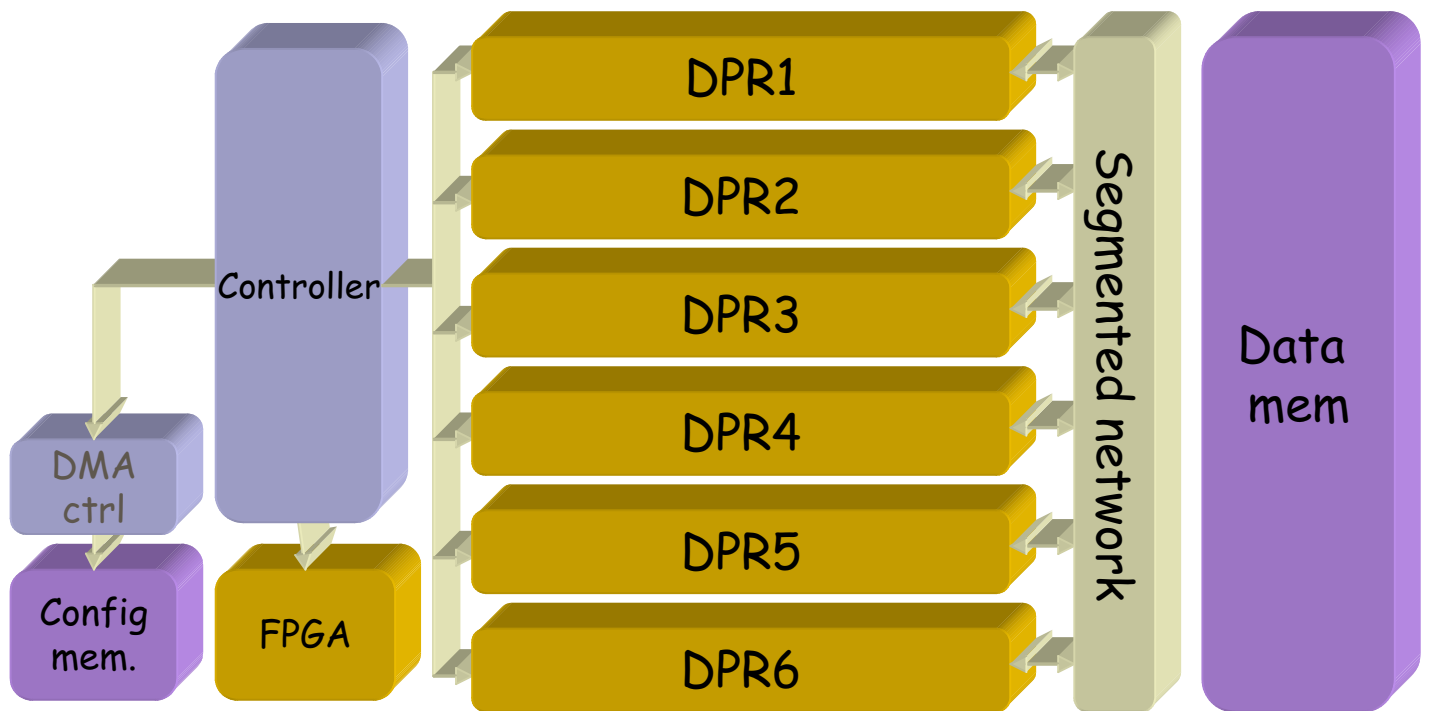
### ☒ Faible consommation

### ☒ Distribution hiérarchique des ressources

- ☒ calcul, stockage, interconnexions, contrôle

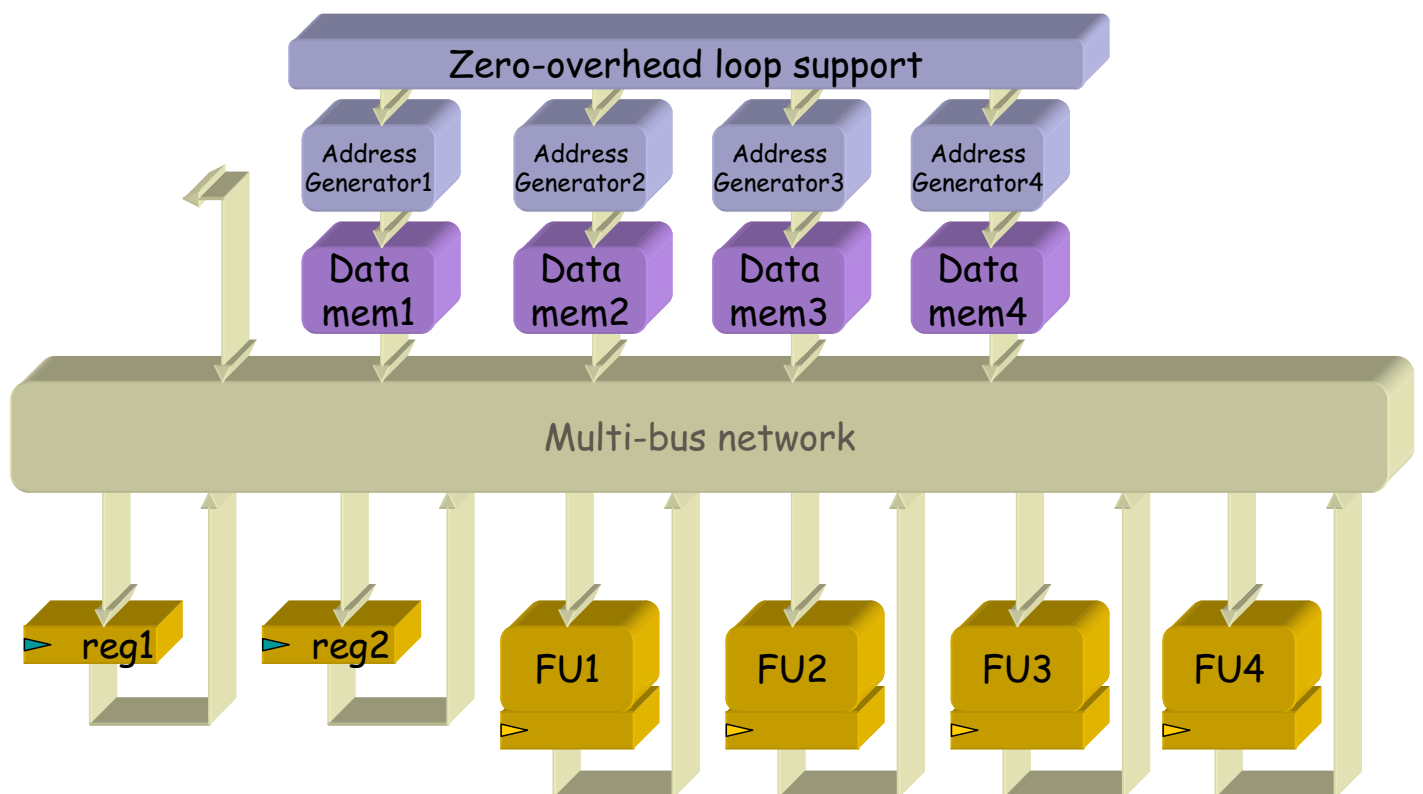


# Dart : une archi reconf



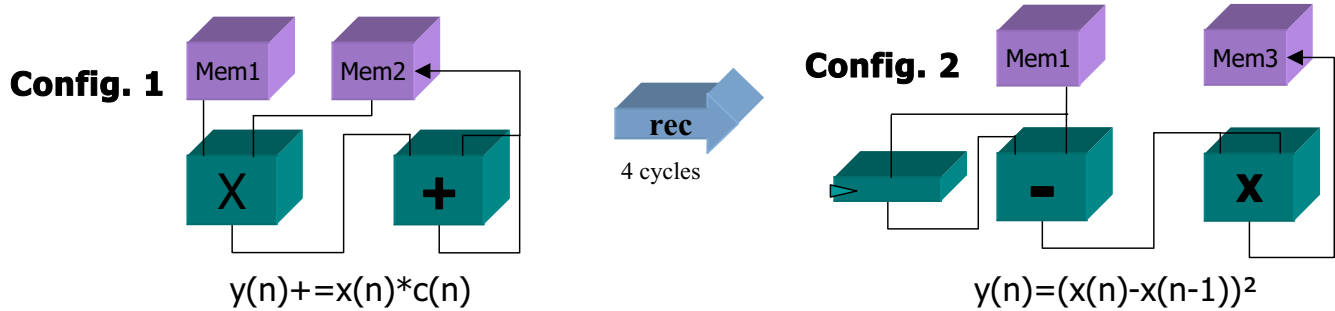
# Dart : une archi reconf

## ⌘ Data Path Reconfigurable

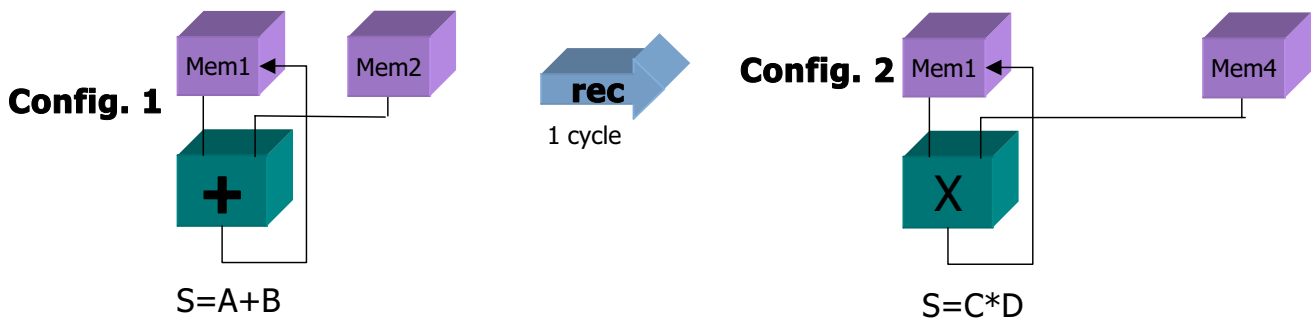


# Dart : une archi reconf

⌘ Reconfiguration HW pour optimiser le chemin de données :

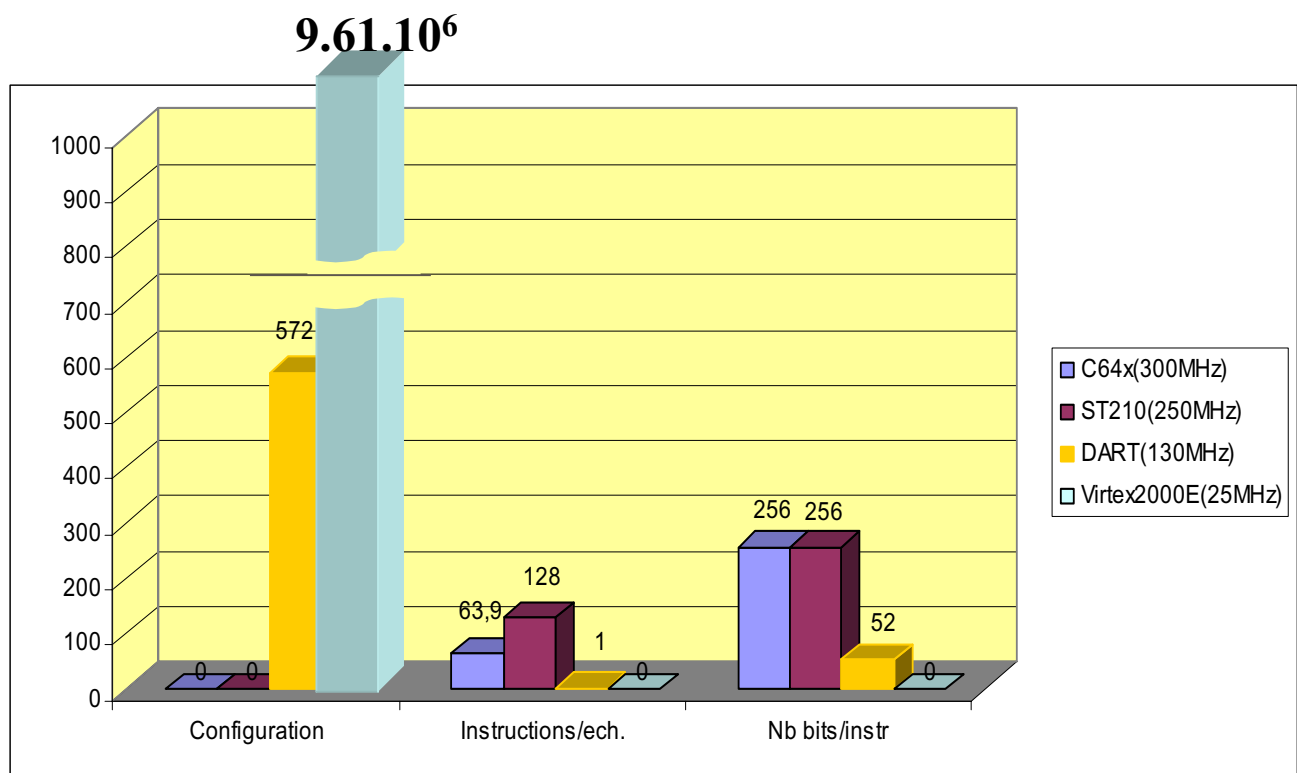


⌘ Reconfiguration SW pour reconfigurer le DPR à chaque cycle :



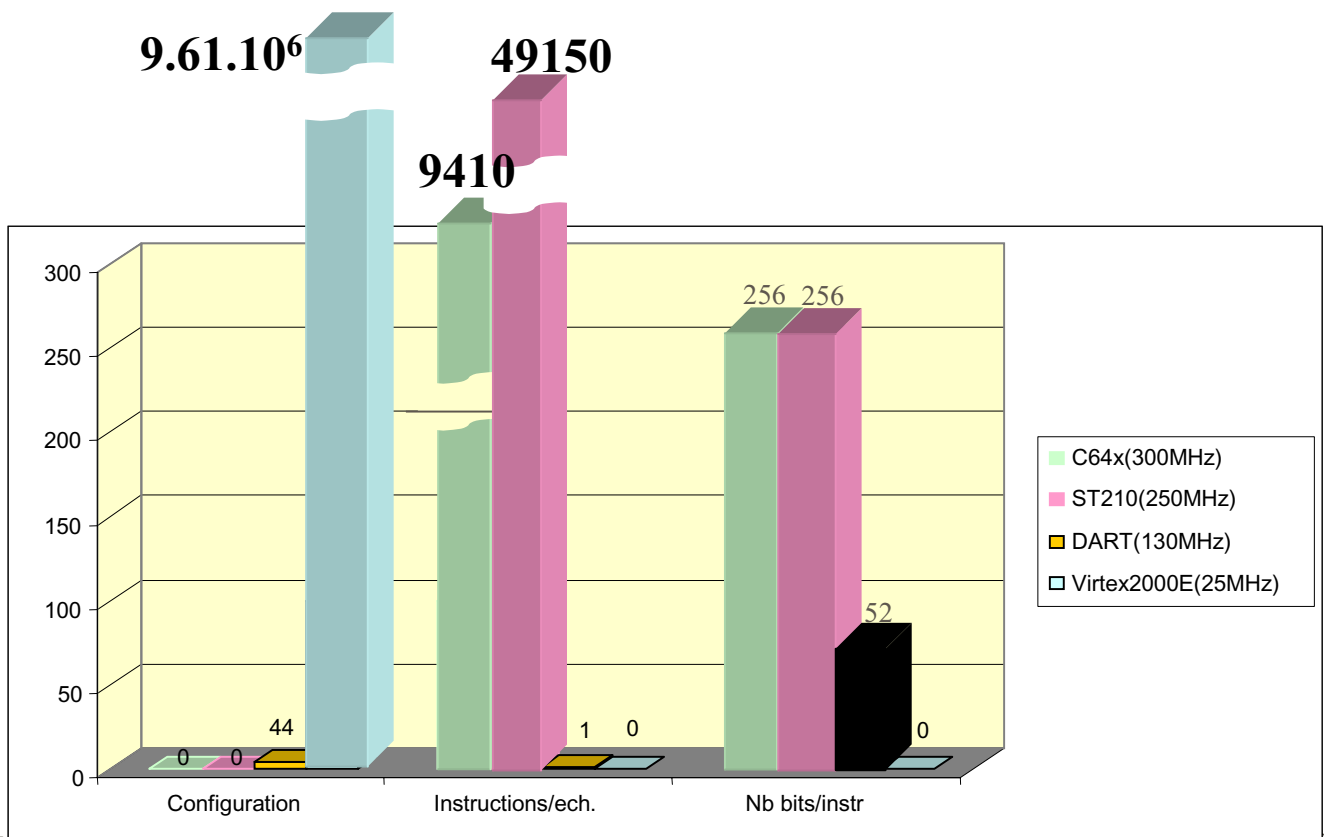
# Dart : une archi reconf

⌘ Résultats : volume des configurations pour le FIR



# Dart : une archi reconf

⌘ Résultats : volume des configurations pour le Rake Receiver

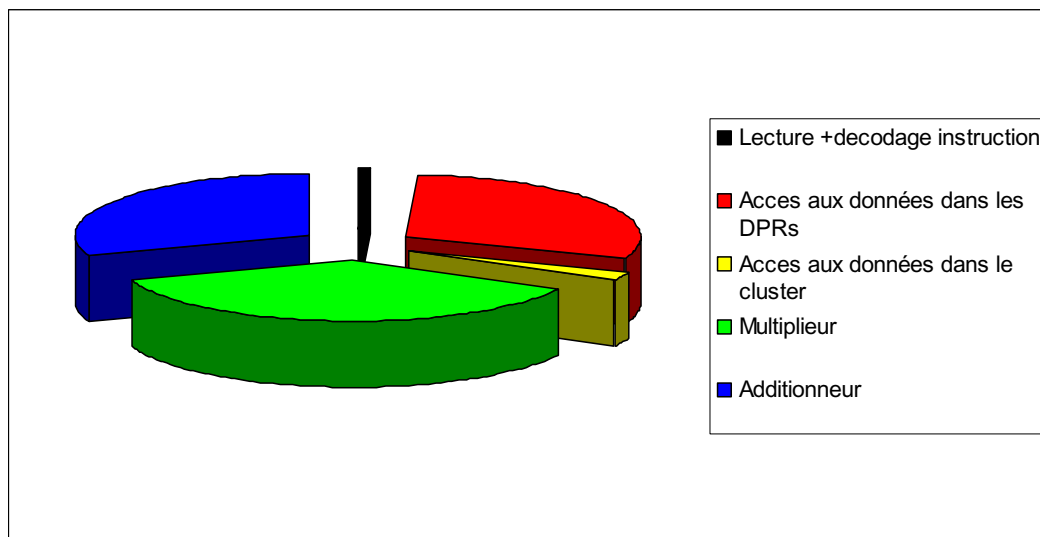


R2D2

Page 87/

# Dart : une archi reconf

⌘ Répartition de la consommation dans DART lors de l'exécution d'un filtrage (1966GMACs)

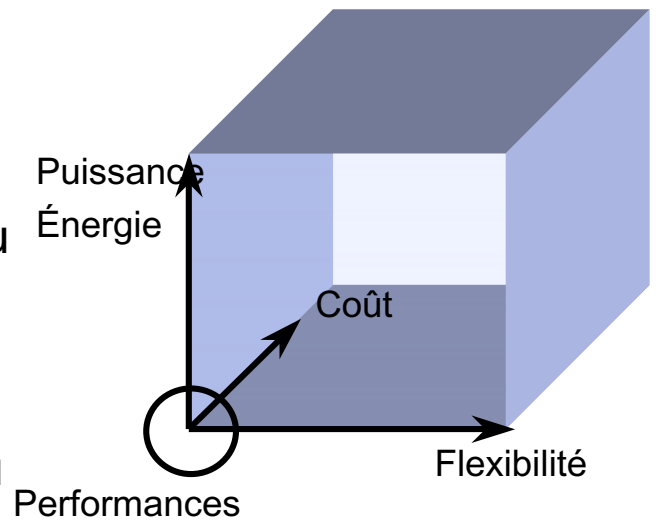


**149mW**

# Conclusions



- ⌘ De nouvelles métriques pour la conception
  - ☑ Puissance ou énergie devient la contrainte dominante
- ⌘ La consommation doit être estimée et optimisée à tous les niveaux, en particulier au niveau système
- ⌘ La diminution de la tension d'alimentation doit se faire en conservant les performances du circuit
- ⌘ La diminution de l'activité de commutation des signaux internes ou externes



Organisation

En collaboration avec



## FTFC'2003

4<sup>ème</sup> journées d'études

Faible Tension Faible Consommation

15-16 mai 2003, Paris, France

Cercle national des armées  
8, place Saint - Augustin, 75008 Paris



IEEE



Circuits and Systems  
Society