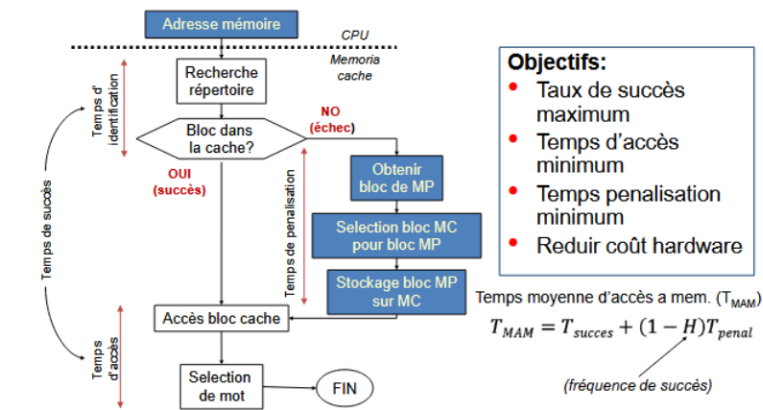


Structure dans le cache

- Le cache est divisé en blocs (ligne) (groupe de 8-64 octets)
- les bloc on une étiquette unique (tag)
- Combien de lignes dans un cache de 64KB avec des lignes de taille 64 Bytes ?
 $KB = 2^4 \cdot 10 = 1024 \text{ (oct)} \rightarrow (64 * 1024) / 64 = 1024$



Fréquence de succès

- dépend de la taille du cache
- Miss penalty devient du Hit time

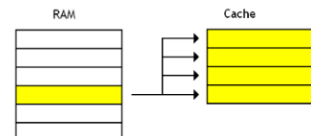
	Rate	Time
Cache Hit	Hit rate = Nombre de hits/Nombre d'accès	Temps d'accès à la cache (hit time)
Cache Miss	1 - (Hit rate)	Temps d'accès à la mémoire principale + temps de chargement cache (miss penalty)

Recherche dans la cache

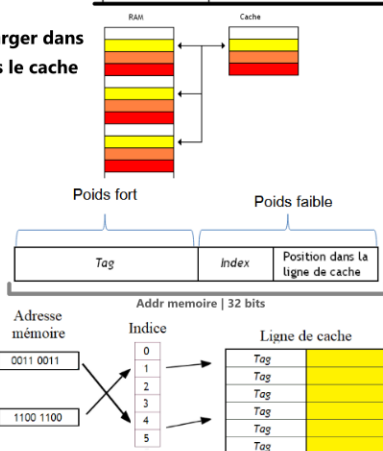
mots = mots dans la memoir principale

bute = trouver la ligne dont le tag correspond à l'Addr de base demander au répertoire

- 1) Complètement associative : mots peuvent être stocker n'importe où dans le cache
 + : taux de hit très élevé
 - : trop lent (séquentiel, toutes les lignes à regarder)



- 2) Associative par ensembles : mots est charger dans des ligne predéfinie (toujours la même) dans le cache
 + : get rapide, car il faut checker 1 ligne
 - : conflit, taux de hit éclater

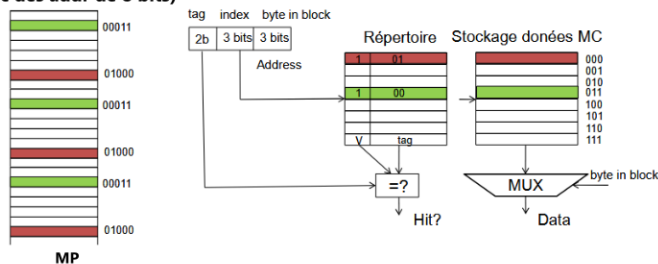


- 3) Associative par voies :
 - Bits faible → index (qu'elle ligne) et position dans la ligne
 - Bits fort → Tag

Chaque bloc tombe sur les même cache

- Mémoire principale (MP) : 256 (oct), 8 (oct) blocs → 32 blocs (Addr 5 (oct))
- Cache (MC) : 64 (oct), 8 blocs (Addr 3 (oct))

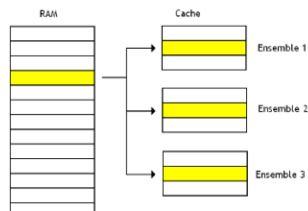
Exemple (avec des addr de 8 bits)



- 4) Associative par voie : Cache composé de plusieurs "cache" (même addr) accessibles en parallèle

- Chaque cache est appelé une voie
- Mots peut être stocké en N positions différentes dans le cache

- Diminution de la fréquence d'échec :
 Associativité *2 = diminution de 20%
 Taille cache *2 = diminution de 69%



Cache (ou antémemoire)

- Composer de transistors et pas de condensateurs
- Pas de contrôle sur son contenu

Localité

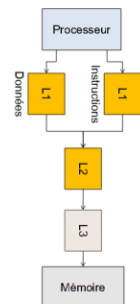
Localité spatiale	Localité temporelle
Le code d'un programme s'exécute toujours à l'intérieur de petites zones répétées de mémoire (des blocs correspondant à des boucles ou/et des sous-programmes). Si une adresse mémoire est utilisée, alors les adresses proches le seront probablement dans un futur proche.	Les blocs s'exécutent en séquences très proches . <i>il y a plus de chances d'accéder à une position de mémoire utilisée il y a 10 cycles qu'à une autre utilisée il y a 10000 cycles</i> Si une adresse mémoire est utilisée, elle le sera probablement à nouveau dans un futur proche

Localité temporelle: Chaque fois que le processeur prend un mot de la mémoire, il copie le mot à la cache.

- Pénalisation dans le premier accès
- Plusieurs utilisations seront plus rapides

Localité spatiale: Au lieu de copier un seul mot, on va en prendre plusieurs en même temps (un bloc)

- Accès à a[0]
- On prend aussi a[1], a[2], a[3]



Niveaux

souvent plusieurs niveaux de cache (instruction et données)

Repertoire

Indique si le mot accédé par le processeur se trouve dans la cache. Si oui → indique l'Addr dans le cache

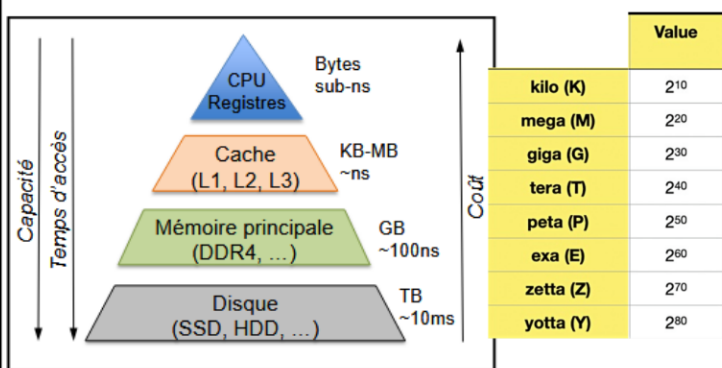
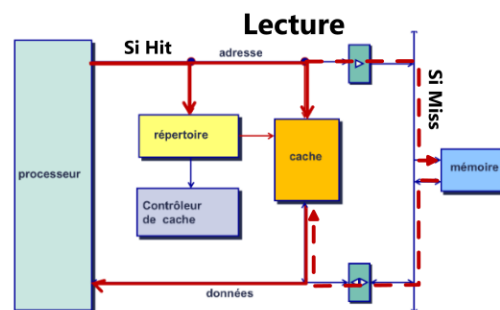
Contrôleur de la cache

Contrôle de la lecture/écriture en mémoire (parmi des buffers, les bloc avec les triangles)

Écriture et Lecture

En gros,

- lecture : quand ça miss ça vas chercher dans la memoire et ça écrit le resulta en même temps que de retourner la valeur
- Ecriture : on écrit dans les 2 de toute manière mais si ça hit, on peut le retourner plus vite



CPU à 2.6 GHz, mémoire RAM avec 10ns de latence. Le CPU veut accéder à une information mémoire.

- Après combien de temps va-t-il recevoir l'information ?

$$10 * 10^{-9} s = 10^{-8} = 10 \mu s$$

- Combien de cycles CPU se sont déroulés pendant ce temps ?

$$10 * 10^{-8} * 2.6 * 10^9 = 26 \text{ cycles}$$

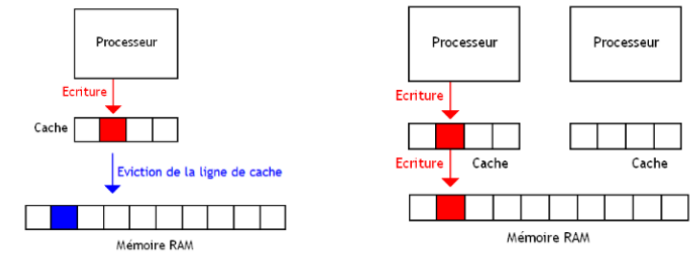


Cache ARM9

- 128KB : 4 voies associative 32KB
- ligne : 64 (oct)

Algo d'écriture

- Write-back
- Write-through



Algo de suppression

- Il faut ajouter l'information manquante à la cache
 - Si la cache est pleine, comment choisir la ligne à remplacer ?

	Fully associative	Direct Mapped	Set associative
Choix			

- Politiques de remplacement

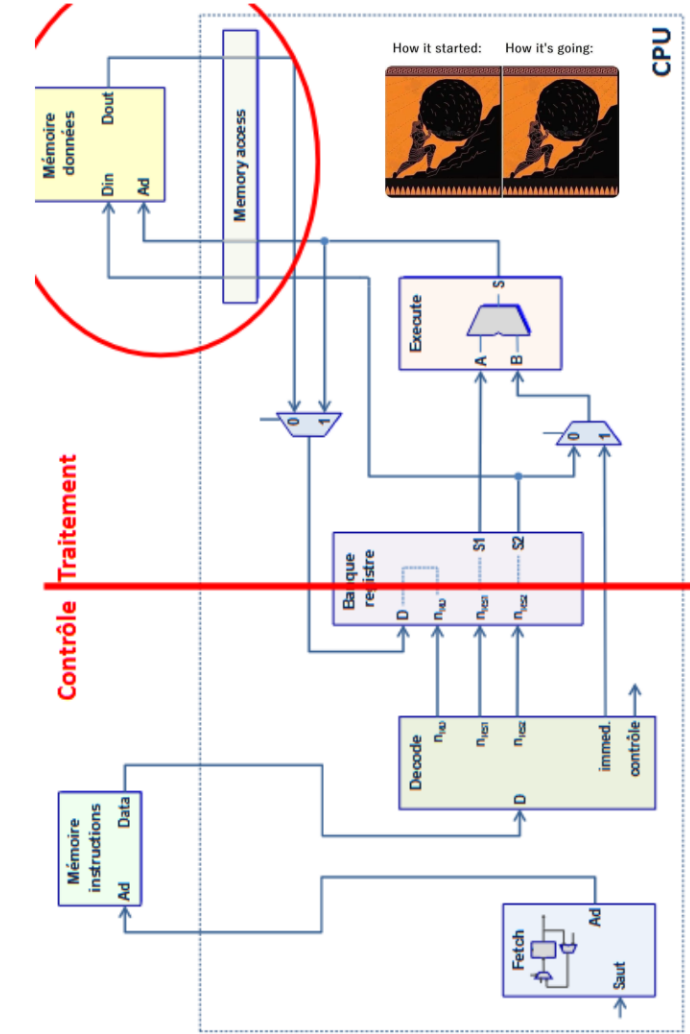
Algo de suppression (cache miss)

- Il faut ajouter l'information manquante à la cache
 - Si la cache est pleine, comment choisir la ligne à remplacer ?

	Fully associative	Direct Mapped	Set associative
Choix	Trop de choix	Pas le choix	Un peu de choix

- Politiques de remplacement

- Aléatoire : on choisit une ligne au hasard
- FIFO : premier entré – premier sorti
- Least Recently Used (LRU) : on supprime une ligne qui n'a pas été utilisée depuis longtemps



Cache Exo methods

Get memorie require 3 (ws)(wait state) et (hr)(hit rate) de 90%

- nombre moyen de waits state = $10\% \times 3 = 0,3$ ws par cycle memoire
- si donné en cache → 0 ws

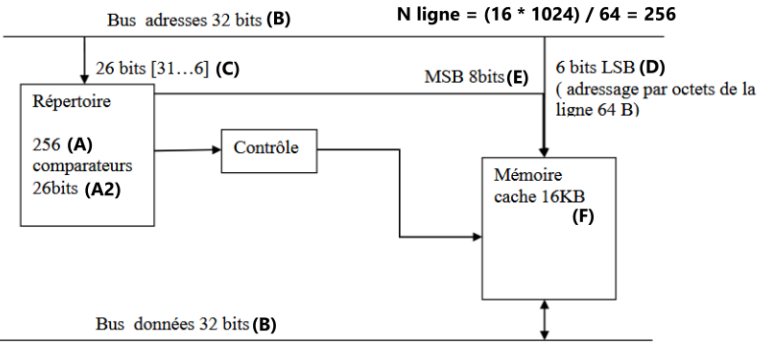
Intruction require 2 cycle sans (ws) et 5 cycle avec | (hr) = 80%

10 instruction → $(10 \times 0,8 \times 2) + (10 \times (1 - 0,8) \times 5) = 16 + 26 = 26$ cycle dont 16 avec cache

Intruction require 2 cycle sans (ws) et 5 cycle avec | (hr) = 80%

10 instruction → $(10 \times 0,8 \times 2) + (10 \times (1 - 0,8) \times 5) = 16 + 26 = 26$ cycle dont 16 avec cache

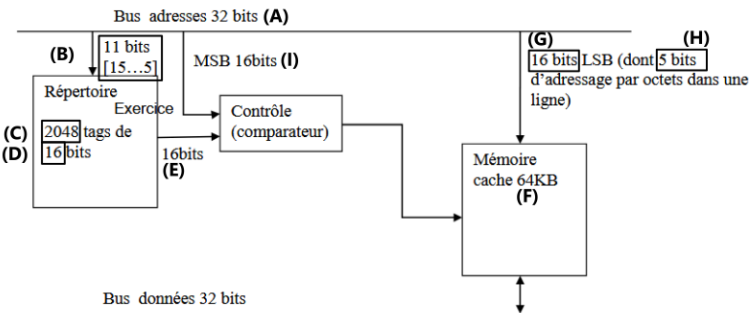
COMPLÈTEMENT ASSOCIATIVE : Bus (Addr et donné) du processeur sont de 32 bits.
Cache de 16KB (Ligne de 64 (oct))



N ligne = $(16 \times 1024) / 64 = 256$

- A) 1 comp par ligne → N ligne
- A2) on a besoin 26 bits à cause de (C)
- B) bus de 32 bits (donné)
- C) 32 - 6 (raison (D)) = 26 bits → [31...6]
- D) offset → $64 = 2^6$ → besoin de 6 bits pour les décrire
- E) N ligne = $256 = 2^8$ → besoin de 8 b pour trouver la ligne
- F) donner dans la donné

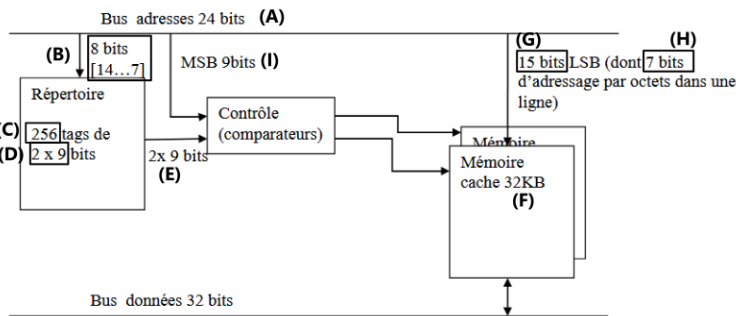
ASSOCIATIVE PAR ENSEMBLES : Bus (Addr et donné) du processeur sont de 32 bits.
Cache de 64KB (Ligne de 32 (oct))



N ligne = $(64 \times 1024) / 32 = 2048$

- A) donner
- B) $(32 - 16) - 5$ (du à (H)) = 11 bits
- C) $(16 - 5) = 11 \rightarrow 2^{11} = 2048 = N$ Ligne → besoin de 16 bits au totale
- D,E) Meme logic que (C)
- F) donner
- G) N ligne = $2048 \rightarrow 11$ bits → $11 + 5$ (du à (H)) = 16
- H) ligne de 32(oct)= $2^5 \rightarrow 5$ bits
- I) $32 - 16 = 16$

ASSOCIATIVE À 2 VOIES : Bus Addr de 24 bits et donné du processeur sont de 32 bits.
Cache de 32KB par voie (Ligne de 128 (oct))



N ligne = $(32 \times 1024) / 128 = 256$

- A) donner
- B) 15 (du à (G)) - 7 = 8 bits
- C) N ligne
- D,E) du à (I) et au fait qu'il y a 2 voie → 2×9
- F) donner
- G) N ligne = $256 \rightarrow 8$ bits → $8 + 7$ (du à (H)) = 15
- H) ligne de 128(oct)= $2^7 \rightarrow 7$ bits
- I) $24 - 15 = 9$

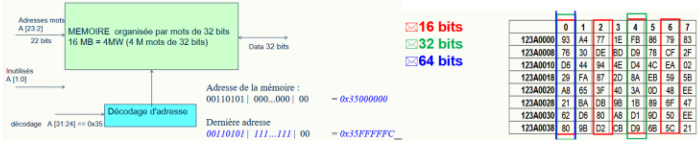
Memoire

- Addr est stockée dans un mot
- N bits d'un mot = taille max de la memoire → mots = 32 bits → 2^32 = 4GB

Décode d'Addr memoire

2^32 = 4MW (4 Mega mots)

Adr = (N * nb_byte)/(taille de mots)



Mémoires volatiles (vives)

- informations **perdues** à la mise hors tension
- lecture et écriture en cours d'utilisation

RAM, SRAM, DRAM

Mémoires non-volatiles (mortes)

- informations **conservées** à la mise hors tension
- lecture en cours d'utilisation
- écriture (« programmation ») durant la fabrication de la mémoire, ou sur la carte (in situ).
- modification du contenu durant le fonctionnement du système, implique un effacement!

EPROM, Flash memory, PROM, FeRAM, EEPROM, ROM, PRAM, E2PROM, MRAM

Memory Map

Exemple :
32 bits d'adressage
SRAM = 32KB = 2^15 = 0111 1111 1111 1111 = 0x00007FFF
Internal Flash = 4MB = 2^22 = 0x007FFFFF
Total mem → 2^32 = 4GB → 0x3FFFFFFF
→ 0xFFFFFFFF - 0x3FFFFFFF = 0xFFC00000
SRAM → 0x00000000 → 0x00007FFF +1

Endianness



Bank de registre

THUMB

R0
R1
R2
R3
R4
R5
R6
R7

R13 (SP)
R14 (LR)
R15 (PC)
CPSR

can you lock the fuck in

Diagram showing register selection and decoding logic.

- SP (Stack pointer) : Stocke la position dans lapile de stockage
- LR (Link Register) : Garde l'adresse de retour(appel de fct)

Pile

- Pile ascendante :
SP init sur l'Addr du bas de la pile

- Pile descendante :
SP init sur l'Addr du haut de la pile

Utilisation Pile

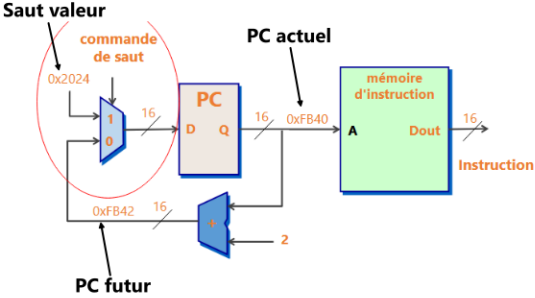
- stockage des Addr de retour pour les appl de fonction et interruption
- Save de registre pour les contexte
- Save de variable locales

Incrémentation

Pré = ++i | Move SP → action
Post = i++ | action → move SP

Fetch

- PC (Programme compteur): incrément chaque instru sauf si saut



Saut inconditionnel

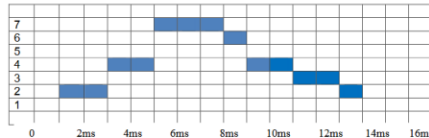
- Adr = PC + extension_16bits(offset(11) x 2) + 4
Exemple : 0xE00F → 1110 0 000 0000 1111
offset(11) x 2 → 000 0001 1110
Extention 16 bits → 0000 0000 0001 1110
+4 → 0000 0000 0010 0010 = v
PC actuel = 0000 0000 0000 1110
v + PC actuel → 0000 0000 0011 0000 → 0x0030

Saut conditionnel

- Adr = PC + extension_16bits(offset 8 x 2) + 4

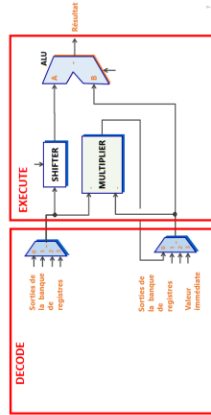
Interruption

- 2) t = 1ms | d = 3ms
- 4) t = 3ms | d = 4ms
- 3) t = 4ms | d = 2ms
- 7) t = 5ms | d = 3ms
- 6) t = 6ms | d = 1ms

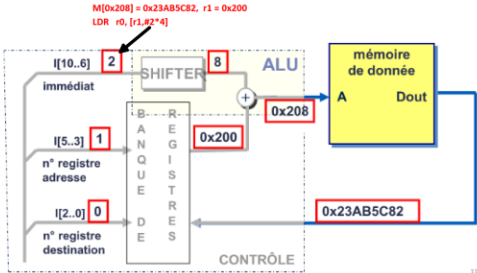


Execute

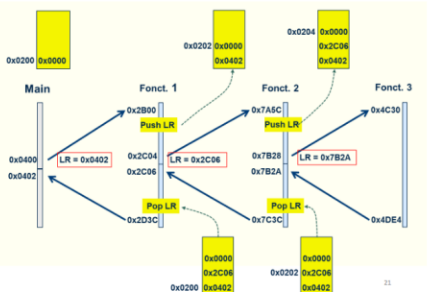
- **Logical shift** : multiplication ou division par une puissance de 2
 - ajout de 0 à droite ou à gauche
- **Arithmetic shift** : multiplication ou division par une puissance de 2 sur des nombres signés
 - ajout de 0 à droite, extension de signe à gauche
- **Rotation** : report des bits retirés d'un coté vers l'autre coté



MEMORY ACCESS

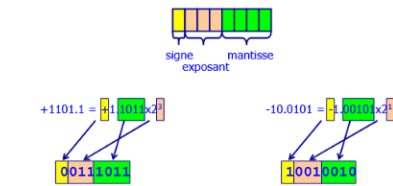
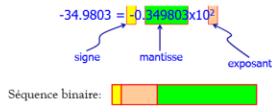


• Un périphérique comprenant 5 registres de 16 bits est placé à l'adresse 0xF230 sur le bus d'un processeur. Le bus d'adresses est un bus de 16 bits. Donnez les adresses des 5 registres (Incrément de 2 car 16 bits) : 0xF230, 0xF232, 0xF234, 0xF236, 0xF238



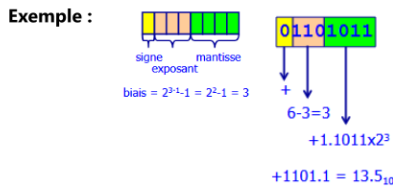
Floa

Structure d'un float



Représentaton avec biais

Champ exposant = exposant + biais
Typiquement "2^(k-1)-1"
avec k=nombre de bits duchamp de l'exposant



Représentaton non normaliser

Pour les nombre très petit proche de 0
+0.0=0 00...00
-0.0=1 00...00
Dans ce cas l'exposant = 1 - biais
Exemple : +0.001 = +0.10 * 2^-2 = 0 000 10



Standard IEEE 754

- normaliser : biais = 127
- Non-normaliser : 1-biais = -126

Decimal → Binaire

0.375 = ?

0.375 x 2 = 0.75
0.75 x 2 = 1.5
0.5 x 2 = 1.0

0.375 = 0.011

0.3 = ?

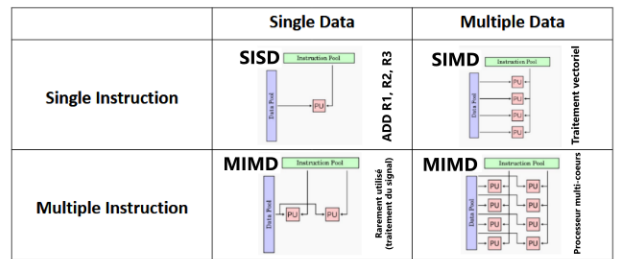
0.3 x 2 = 0.6
0.6 x 2 = 1.2
0.2 x 2 = 0.4
0.4 x 2 = 0.8
0.8 x 2 = 1.6
0.6 x 2 = 1.2

0.3 = 0.01001[1001] = 0.0[1001]

exposant	1	2	3	4	5	6	7	infini
champ exposant	1000	1001	1010	1011	1100	1101	1110	1111
exposant	non normalisé	-6	-5	-4	-3	-2	-1	0
champ exposant	0000	0001	0010	0011	0100	0101	0110	0111



Pipline



Aléas

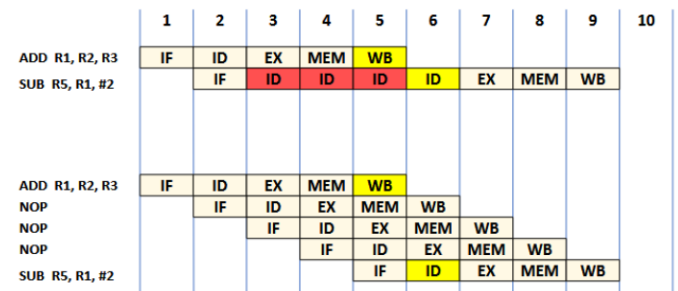
Dépendances

- Structurel : Accès aux meme ressources
- Données : Utilisation des même registres
- Contrôle : Décision de branchement

Dépendances de données

- RAW (read after write) = ADD R1, R2, R3
SUB R5, R1, #2
- WAR (write after read) = ADD R1, R2, R3
SUB R2, R4, #2
- WAW (write after write) = ADD R1, R2, R3
AND R5, R1, R2
SUB R1, R4, #2

Stop pipeline



Forwarding

Cours ARO 2025

Exemple de pipeline avec bypass

r4 = 0x12, r3 = 0x28B2,
M[0x28BA] = 0xABCD



Pipline MATH

Base

- T_e = temps de passage à chaque étape
- T_p = n*T_e = temps pour n étapes

Calcul de performances

- Latence : Temps entre le debut et la fin de l'exec
- Débit : nombre d'op executer par unités de temps

$$D = \frac{m}{(n+m-1) * T_e} \approx \frac{1}{T_e}$$

-Accélération : n fois plus rapide que le séquentiel

$$A = \frac{T_{seq}}{T_{pip}} = \frac{m+n*T_e}{(n+m-1)*T_e} = \frac{m+n}{n+m-1} \sim n \text{ (pour m très grand)}$$

- IPC (nombre d'instru par cycle)

Exemple : stop 3 cycle pour 20% des instru

$$IPC = 1 / (0.8 \times 1 + 0.2 \times 4) = 0.625$$

