

ORGANIZAÇÃO E ARQUITETURA DE COMPUTADORES

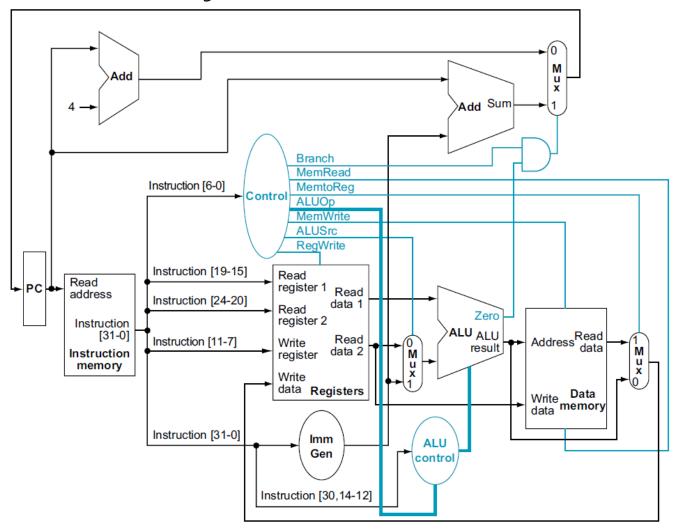
Arquitetura RISC-V Pipelined

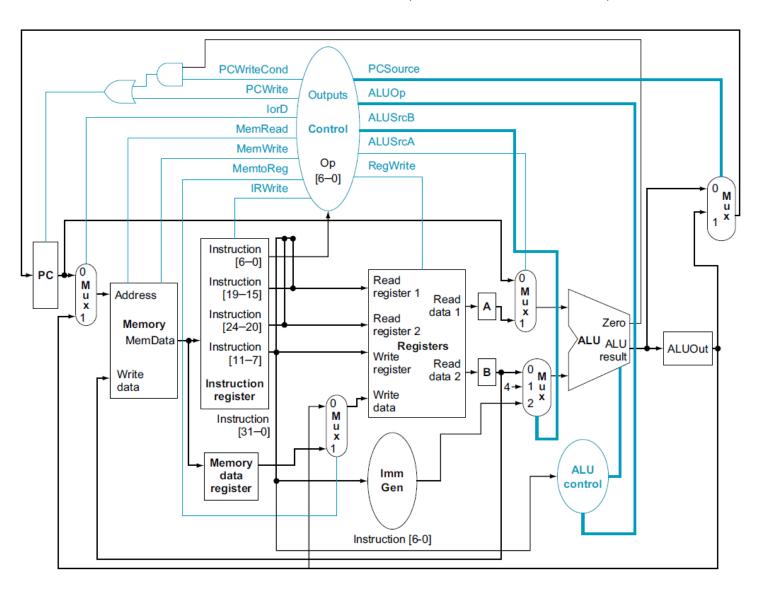
Prof^a. Fabiana F F Peres

Apoio: Camile Bordini

Monociclo (revisão)

• Executa toda instrução em um ciclo de clock;





(características)

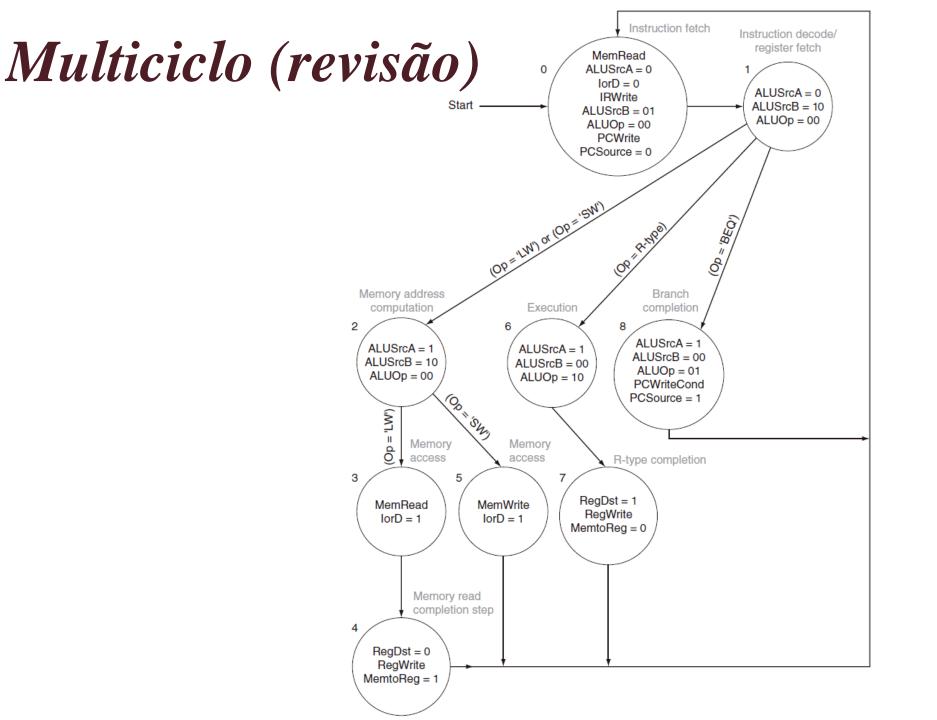
- •Cada um dos passos, tomará um ciclo de clock;
- •Permite que as unidades funcionais possam ser usadas mais de uma vez pela mesma instrução contanto que seja em ciclos de clock diferentes;
- •Reduz o hardware;
- •Permite que as instruções tomem diferentes números de ciclos de clock;

(passos para execução das instruções)

- 1. Buscar instrução da memória e incrementar PC;
- 2. Decodificar a instrução, ler registradores e calcular endereço de desvio condicional;
- 3. Depende da instrução que está sendo executada:
 - Load ou Store: calcular endereço da memória;
 - Formato R: realizar a operação logica ou aritmética requisitada;
 - Beq: realizar a comparação entre o conteudo dos dois registradores e realiza o desvio se a cond. for aceita;
 - Jump: atualizar o PC com o endereço de desvio;

(passos para execução das instruções)

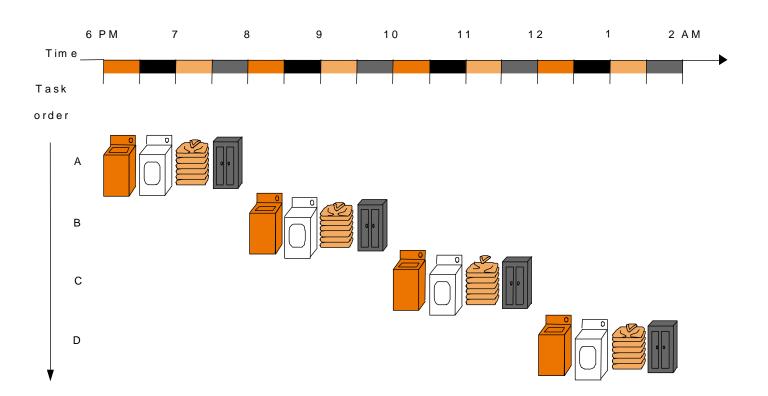
- 4. Depende da instrução que está sendo executada:
 - Load: Lê o dado da memória;
 - Store: Escreve o dado na memória;
 - Formato R: Escreve o resultado da operação no conjunto de registradores;
- 5. Depende da instrução que está sendo executada:
 - Load: Escreve o dado lido no conjunto de registradores;



Multiciclo

(analogia)

- Instrução: "lavar roupa";
- Etapas: lavar, secar, passar e guardar;
- Execução das tarefas **de forma sequencial** (<u>leva 8h</u>):



Pipeline

(definição)

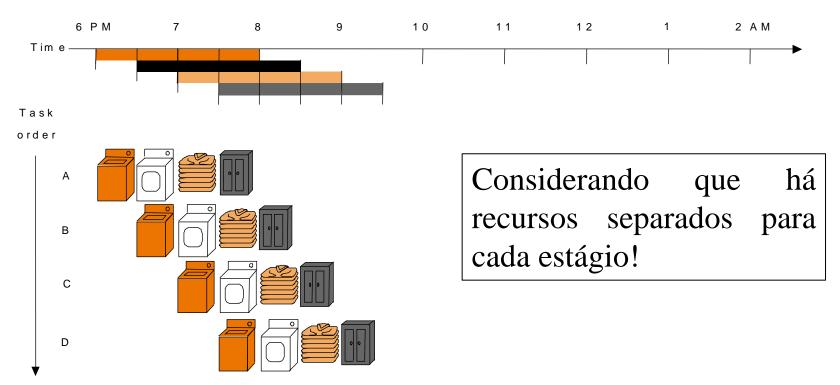
•*Pipelining* é uma técnica de implementação onde várias instruções são executadas simultaneamente

•Atualmente o pipeline é quase universal

Pipeline

(analogia)

- Instrução: "lavar roupa";
- **Etapas**: lavar, secar, passar e guardar;
- Execução das tarefas **de forma paralela** (<u>leva 3,5h</u>):



Tempo das instruções

Monociclo

Tempo do ciclo = tempo total da instrução mais longa

Multiciclo

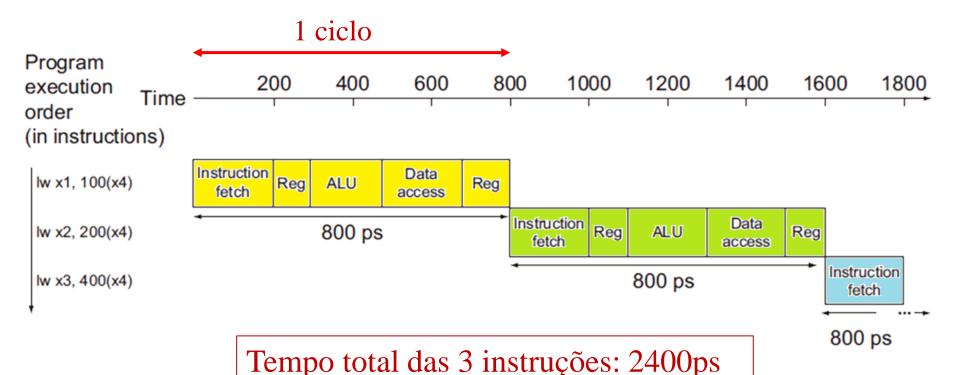
- Tempo do ciclo = tempo total da etapa mais longa
- Instruções executam <u>apenas etapas necessárias</u> (nº de ciclos entre 3 e 5)

Pipeline

- Tempo do ciclo = tempo total da etapa mais longa
- Instruções executam todas as etapas, até mesmo as desnecessárias (nº de ciclos fixo: 5)

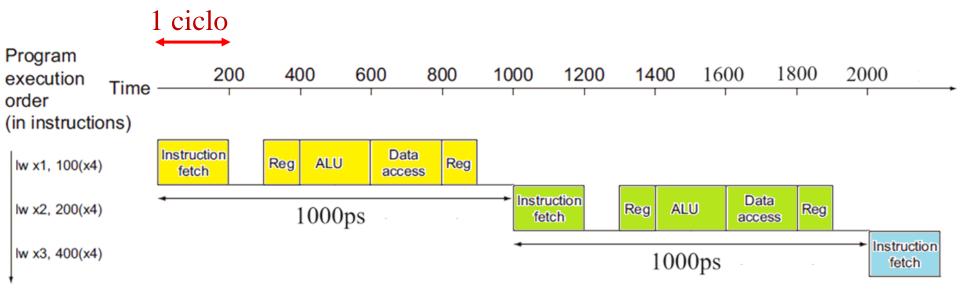
Tempo das instruções - Monociclo

Instruction class	Instruction fetch	Register read	ALU operation	Data access	Register write	Total time
Load word (lw)	200 ps	100 ps	200 ps	200 ps	100 ps	800 ps
Store word (sw)	200 ps	100 ps	200 ps	200 ps		700 ps
R-format (add, sub, and, or)	200 ps	100 ps	200 ps		100 ps	600 ps
Branch (beq)	200 ps	100 ps	200 ps			500 ps



Tempo das instruções - Multiciclo

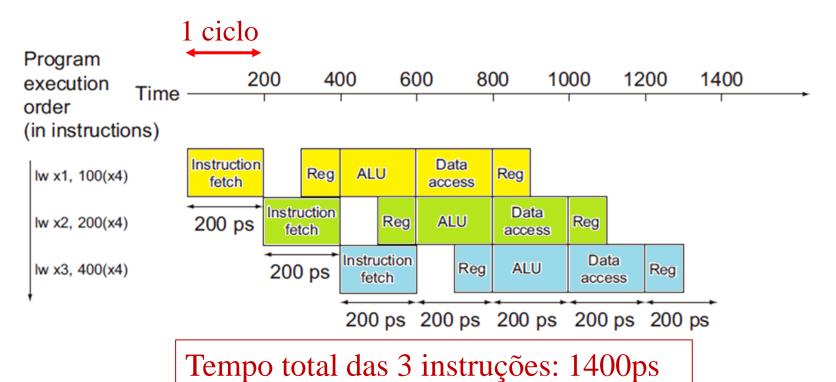
Instruction class	Instruction fetch	Register read	ALU operation	Data access	Register write	Total time
Load word (lw)	200 ps	100 ps	200 ps	200 ps	100 ps	800 ps
Store word (sw)	200 ps	100 ps	200 ps	200 ps		700 ps
R-format (add, sub, and, or)	200 ps	100 ps	200 ps		100 ps	600 ps
Branch (beq)	200 ps	100 ps	200 ps			500 ps

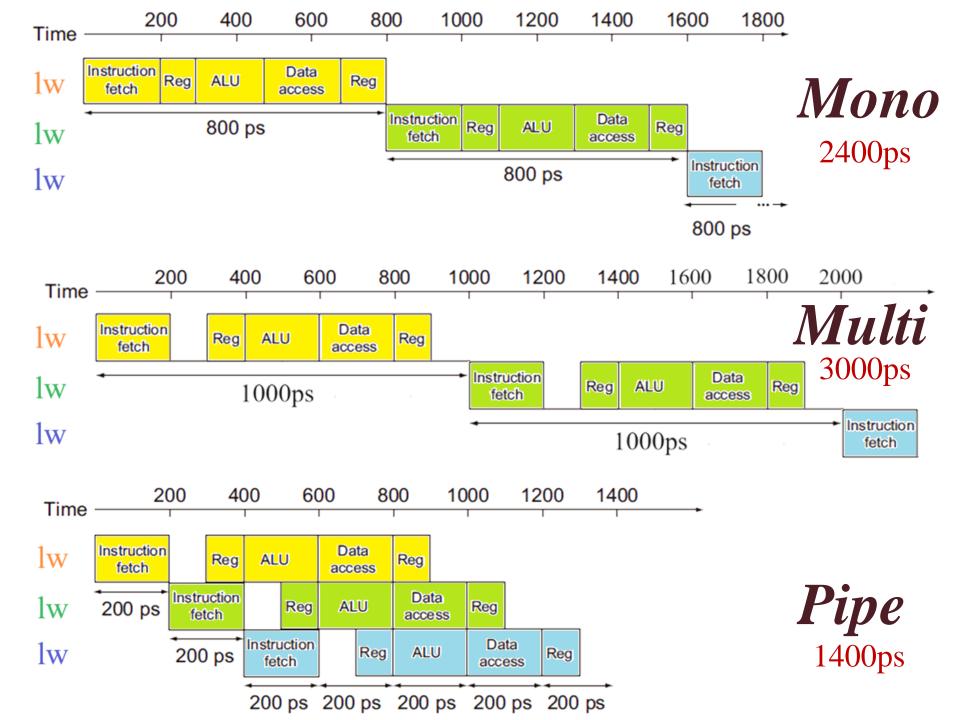


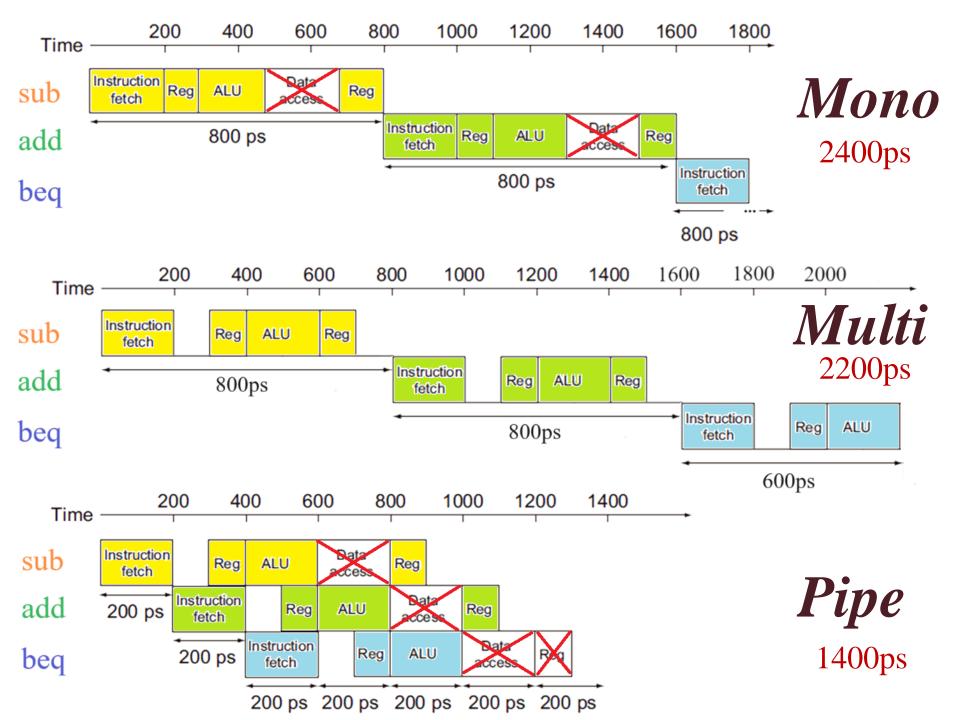
Tempo total das 3 instruções: 3000ps

Tempo das instruções - Pipeline

Instruction class	Instruction fetch	Register read	ALU operation	Data access	Register write	Total time
Load word (lw)	200 ps	100 ps	200 ps	200 ps	100 ps	800 ps
Store word (sw)	200 ps	100 ps	200 ps	200 ps		700 ps
R-format (add, sub, and, or)	200 ps	100 ps	200 ps		100 ps	600 ps
Branch (beq)	200 ps	100 ps	200 ps			500 ps







Pipeline

(definição)

- *Questão*: o tempo total de execução de uma instrução é menor no pipeline?
 - O pipelining não reduz o tempo total da execução (latência) de uma instrução isolada, e sim o throughput (vazão) da carga de trabalho

"Para um pipeline de **5 estágios**, o fluxo de instruções será **5 vezes maior** que em uma implementação sequencial"

•Ou seja, se todas as etapas levassem aproximadamente a mesma quantidade de tempo e houver trabalho em todas as etapas, então a aceleração (*speed-up*) devido ao pipeline é igual ao número de estágios existentes

("Mundo ideal")

•Essa reflexão sobre o *speed-up* do pipeline pode ser vista através de uma fórmula (assumindo as condições ideais anteriores)

$$Tempo \, entre \, instruções_{compipeline} = \frac{Tempo \, entre \, instruções_{sempipeline}}{N\'umero \, de \, est\'agios \, do \, pipe}$$

•Como ficaria no caso do exemplo anterior com as instruções load's?

Tempo entre instruções
$$=\frac{800}{5}$$
 = ciclo de clock de 160ps

•*Para pensar*: uma situação em que o nº de <u>tarefas</u> não é grande em relação ao nº de <u>estágios</u>, como isso afetaria o *pipeline*, para melhor ou pior?

•Lembrando: pipeline é uma forma de implementação (não está vinculado à arquitetura X ou Y)

- •*Para pensar*: uma situação em que o nº de <u>tarefas</u> não é grande em relação ao nº de <u>estágios</u>, como isso afetaria o *pipeline*, para melhor ou pior?
 - -Para pior, quando o *pipeline* não está completamente cheio, isso afeta o desempenho (caso do 1º exemplo com só 3 load's)
 - Mas programas reais executam bilhões de instruções
 - -Com os estágios completamente cheios na maior parte do tempo e um grande número de instruções, o aumento no rendimento deve se aproximar do nº de estágios.

- •No entanto, o que faz o *speed-up* do *pipeline* **não** aumentar proporcionalmente ao número de estágios?
 - •...os estágios não serem perfeitamente equilibrados
 - •...há também conflitos existentes entre as instruções (veremos mais adiante)
 - •...os componentes de interligação necessários p/ implementar um *pipeline* e corrigir esses conflitos introduzem uma latência adicional.

Pipeline + RISC-V

Pipeline RISC-V

•Como o conjunto de instruções (ISA) foi projetado para a execução em pipelining?

- a. Instruções de tamanho variável **vs.** todas instruções de mesmo tamanho?
 - Em RISC-V todas as instruções têm o mesmo tamanho: facilidade em buscá-las no 1º estágio e decodificar no 2º
 - Exemplo: nas primeiras implementações do x86 (instruções variavam de 1 à 15 bytes), o pipeline é bem mais desafiador.

Pipeline RISC-V

•Como o conjunto de instruções (ISA) foi projetado para a execução em pipelining?

- b. Operandos que residem em memória em qualquer instrução vs. operandos residindo em memória apenas através de load e store?
 - Na RISC-V só aparecem em Loads e Stores
 - X86: operandos residindo em memória são possíveis em instruções de operações, isso requereria em pipeline: estágio de endereço -> estágio da memória -> estágio de execução

Pipeline RISC-V

•Como o conjunto de instruções (ISA) foi projetado para a execução em pipelining?

- c. Formato das instruções irregular vs. regular?
 - RISC-V possui apenas alguns formatos de instrução, com registradores (de origem e destino) localizados no mesmo lugar em cada instrução

Passos das instruções em RISC-V

- 1. Busca da **instrução** e incremento do PC
- 2. Acesso ao banco de registradores e **decodificação** da instrução
- **3.** Execução da operação ou cálculo do endereço de memória
- 4. Acesso à **memória** de dados (para leitura ou escrita), se necessário.
- **5. Escreve** o resultado dentro de um **registrador**, se necessário

- 1. Busca da **instrução** e incremento do PC → IF: Intruction Fetch
- 3. Execução da operação ou cálculo do endereço de memória

 EX: Execução
- 4. Acesso à **memória** de dados (para leitura ou escrita), se necessário.

 MEM: Memória
- **5.** Escreve o resultado dentro de um registrador, se necessário

 WB: Write Back

(via de dados)

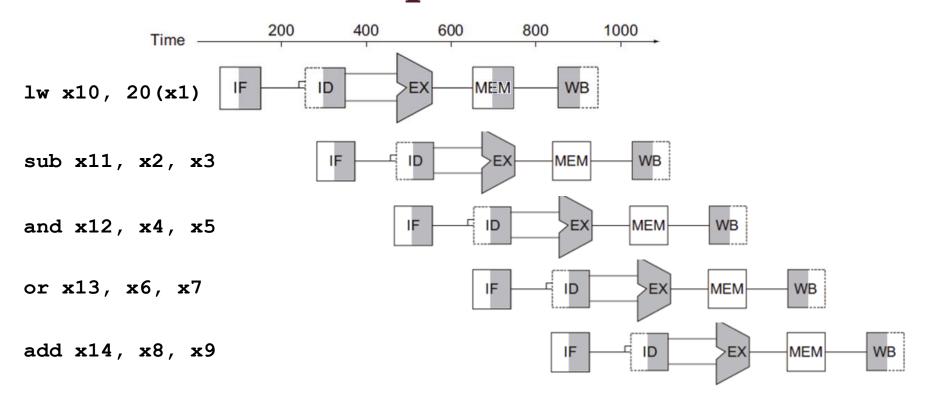
- •Durante os ciclos de clock, cinco instruções estarão em execução;
- •Divisão da via de dados em cinco partes, onde cada uma delas terá o nome do estágio correspondente:
 - •IF: busca a instrução;
 - •ID: decodifica instrução e faz leitura de registradores;
 - •EX: executa operação ou calcula endereço;
 - •MEM: acessa dados da memória (se necessário);
 - •WB: escreve resultado em registrador (se necessário).

Pipeline

(execução de instruções em paralelo)

Program	Time (in clock cycles)						
execution order (in instructions)	CC 1	CC 2	CC 3	CC 4	CC 5	CC 6	
lw \$10, \$20(\$1)	IF	ID	EX	MEM	WB		
sub \$11, \$2, \$3		IF	ID	EX	MEM	WB	

Pipeline



- •Sombreamento indicando que o elemento é utilizado pela instrução
 - -Sombreamento à direita no banco de registradores ou memória: o elemento é **lido** neste estágio
 - –Sombreamento à esquerda: elemento foi escrito neste estágio ³¹

Conflitos

Conflitos

- •Há situações no pipeline que a próxima instrução não pode ser executada no ciclo de clock seguinte
- •Esses eventos são chamados de **conflitos** (*hazards*)
- •Há 3 tipos de conflitos:
 - 1. Estruturais
 - 2. De dados
 - 3. De Controle

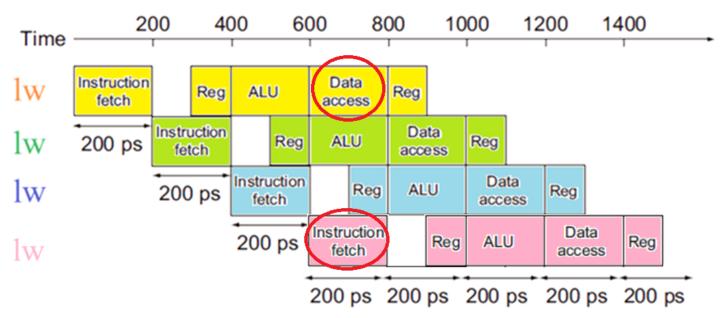
Conflitos estruturais

1. Conflitos Estruturais

- •Situação onde o hardware não suporta a combinação de instruções que se deseja executar num mesmo ciclo de clock
 - Analogia: ocorreria um conflito estrutural se a máquina de lavar e secar fosse a mesma e quiséssemos lavar um conjunto de roupas e secar um outro conjunto
- •O conjunto de instruções **RISC-V** foi projetado para ser *pipelined*!
 - -Isso torna fácil aos projetistas evitar riscos estruturais

1. Conflitos Estruturais

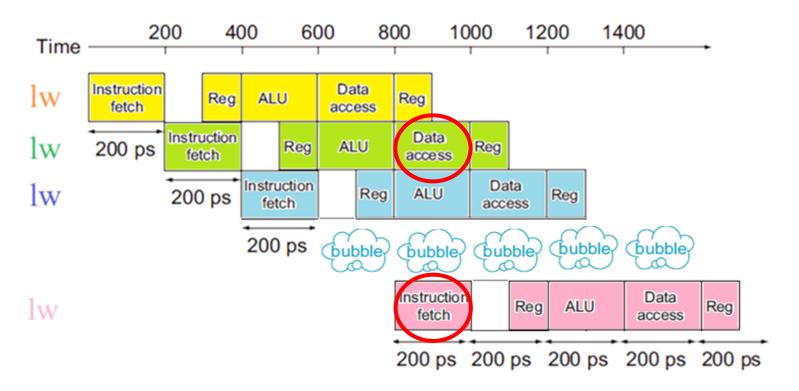
•Suponha: se houvesse apenas <u>uma memória</u> (ao invés de duas). No *pipeline* abaixo, se tivesse uma 4ª instrução **load** o que ocorreria?



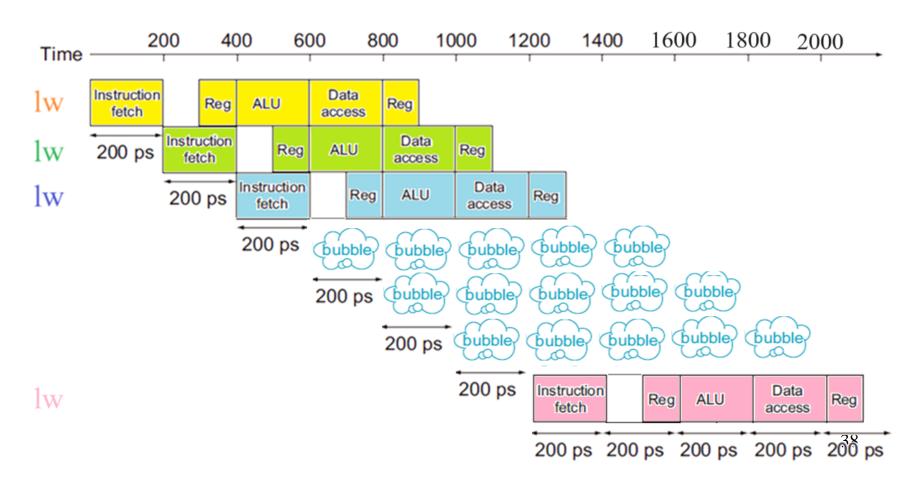
•No 4° ciclo: a primeira instrução tenta acessar dados da memória e a quarta instrução estaria buscando uma instrução da mesma memória!

•Solução 1: inserir uma "bolha". Problema resolvido?

-Passamos a ter problema com a **load** seguinte!



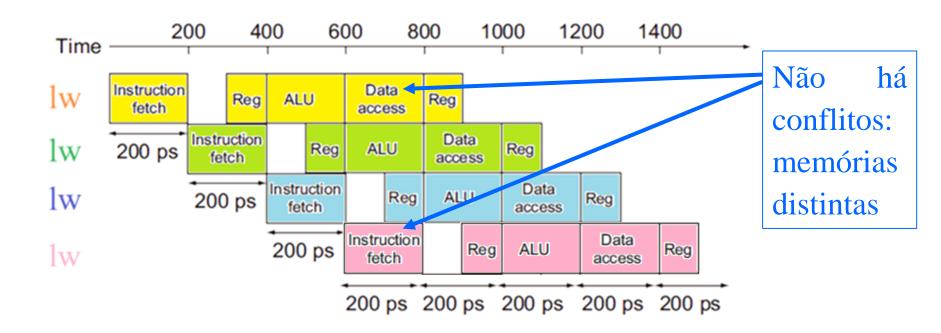
•Ou seja, a cada 3 **load's** seguidas, seria necessário 3 bolhas para se inserir uma 4ª instrução



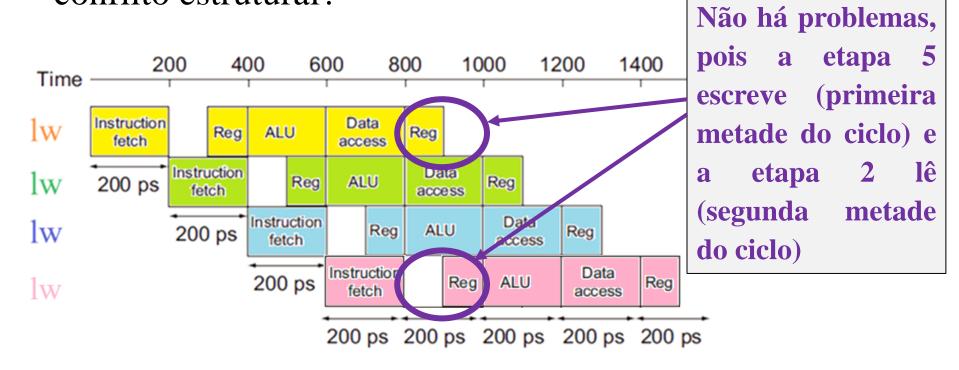
•Solução 2: replicar unidades funcionais para eliminar o conflito e permitir execução simultânea.

- •Como decider se vale a pena replicar ou não?
 - -Mais unidades funcionais encarecem o projeto
 - -Mais unidades funcionais podem alterar o caminho crítico do projeto, e portanto, o tempo do ciclo
- •Qual a frequência da ocorrência dos conflitos estruturais daquele tipo?

•*Exemplo*: com duas portas de leitura na memória: uma para <u>instruções</u> e outra para <u>dados</u> (= Monociclo). Desta forma são consideradas unidades funcionais distintas.



•Questão: e no caso do banco de registradores sendo utilizado tanto pela 1^a instrução quanto pela 4^a, há conflito estrutural?



Conflitos de dados

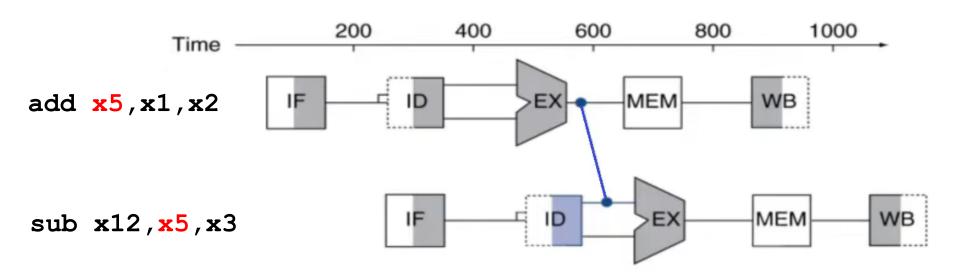
- •Conflitos de dados ocorrem quando o pipeline precisa ser paralisado porque uma etapa precisa esperar outra completar
 - -Analogia: você encontra uma meia sem par na hora de dobrar, e a outra ainda está sendo lavada. Ela precisa passar pelas etapas de secar e passar antes de ter acesso.
- •No *pipeline*: dependência de uma instrução por outra que ainda está no *pipeline*. Exemplos?

•Exemplo: uma instrução add seguida por uma sub que usa a soma anterior: add x19, x0, x1

sub
$$x2$$
, $x19$, $x3$

- •A add não escreve o resultado enquanto não chega no 5° estágio. Soluções?
 - -Desperdiçar 3 ciclos de clock no pipeline (stall no pipe)
 - -Não esperar a conclusão da add para resolver o conflito: assim que a ULA calcular a soma, fornecê-la como entrada para a sub. -> Necessário adicionar hardware extra
 - -Confiar nos compiladores para remover esses perigos

•Solução 1: encaminhar o dado, se possível (forwarding)



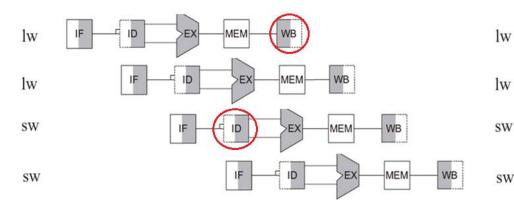
- •Saída da ULA é encaminhada para entrada da propria ULA, sem precisar salvar em registrador antes
- •Pois a **sub** precisa do dado atualizado no início da 3ª etapa

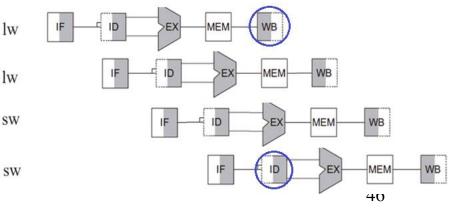
•Solução 2: solução em software pelo compilador: reordenar o Código para evitar bolhas



```
lw x2, 0(x1)
lw x4, 4(x1)
sw x2, 0(x1)
sw x3, 4(x1)
```

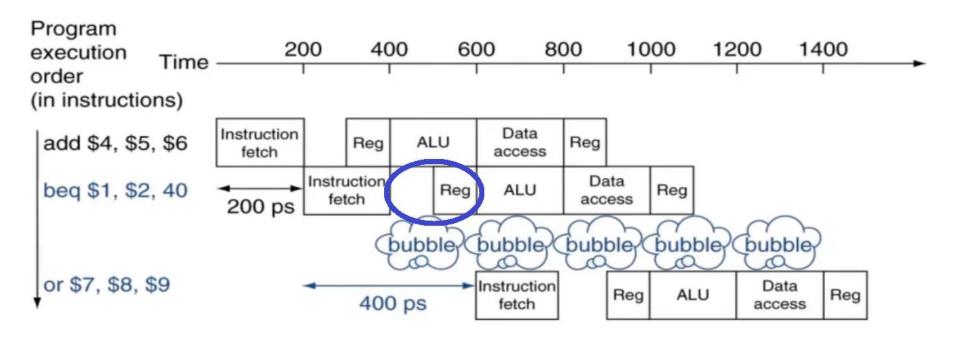
Código reordenado:



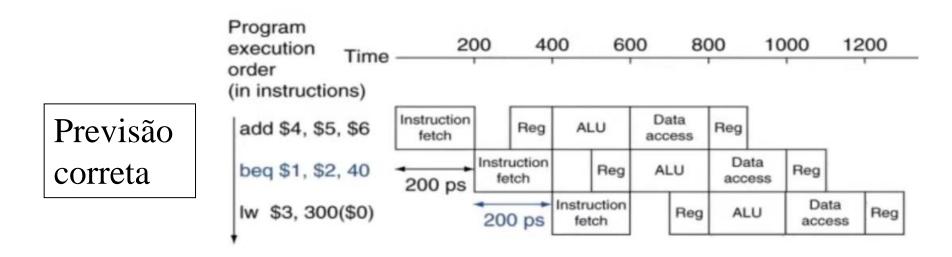


- •Necessidade de tomar a decisão com base nos resultados de uma instrução, enquanto outras estão sendo executadas ainda no pipeline (branches e jumps);
 - -Analogia: a lavanderia precisa limpar os uniformes de um time de futebol, mas como está muito suja a roupa, é preciso determinar o tipo do sabão e a temperatura da água ideais. Depois da primeira roupa seca, é possível avaliar o resultado e mudar as configurações e lavar novamente (*stall*)
 - -Prever/supor que será necessário produtos adicionais de fato (previsão de desvio)
 - -Ou colocar roupas que não sejam de futebol até que a decisão seja tomada (*delayed branch*) 48

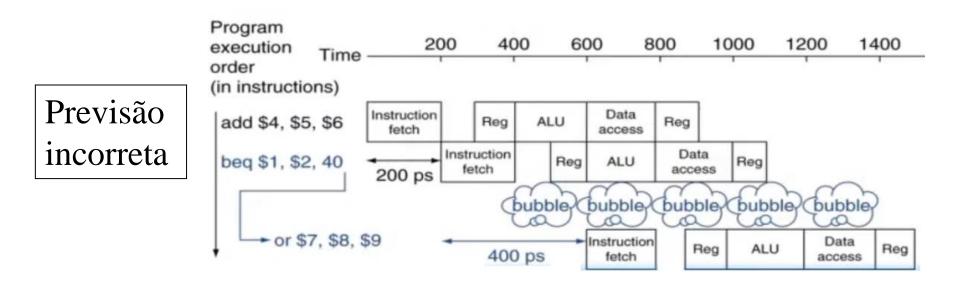
- •Com a adição de hardware: cálculos de desvio passam a ser executadas no 2º estágio (ID)
 - -Mesmo assim o problema permanece
- •Solução 1: necessário 1 bolha (stall)



•Solução 2: previsão de desvio — ex: supor que o desvio nunca será tomado

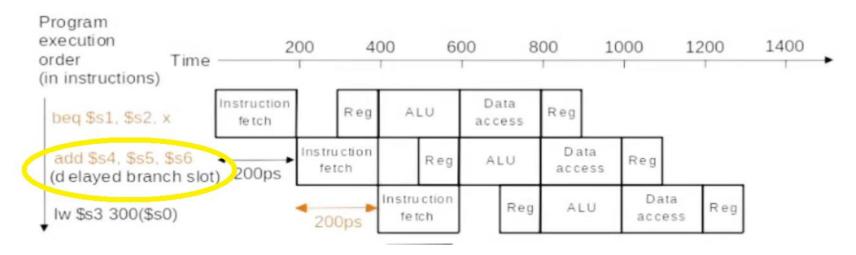


•Se a previsão for correta, não é necessário bolhas



•Se a previsão for incorreta — necessário bolha para anular a instrução **lw** anterior

•Solução 3: delayed branch — colocar uma instrução (independente) que era anterior ao branch para executar logo depois do branch, preenchendo a bolha (*delay slot*) — Feito pelo compilador



•Ex: **add** era uma instrução anterior à **beq**, mas foi colocada após, fazendo a vez de uma bolha

Referências

• PATTERSON, David A; HENNESSY, John

L; Computer Organization and Design – The hardware/software interface RISC-V edition; Elsevier – Morgan Kaufmann/Amsterdam.

• PATTERSON, David; Waterman, Andrew;

The RISC-V reader: an open architecture atlas; First edition. Berkeley, California: Strawberry Canyon LLC, 2017.