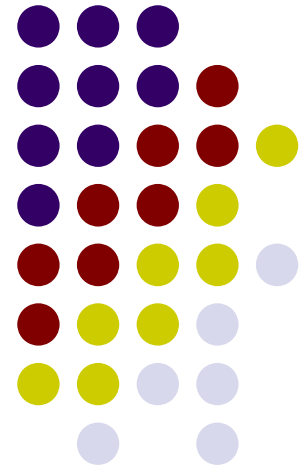


# Transações em Bancos de Dados Relacionais e NO-SQL

Prof. Antonio Guardado

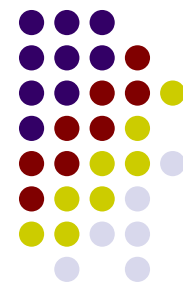


# Agenda

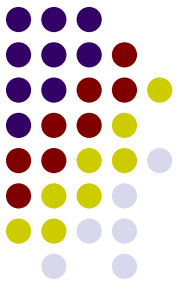


- **1- Gerenciamento de Transações e Propriedades Transacionais nos BDs Relacionais – ACID**
- **2 - Serialização**
- **3- Propriedades Transacionais nos BDs NO-SQL**
- **4 – Propriedades BASE**
- **5 – Teorema CAP**
- **6 – Sistemas CAP**
- **7 – Visão Geral CAP**

# Objetivos

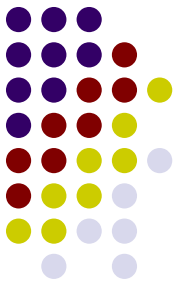


- Entender os problemas das transações distribuídas em BDs Relacionais x Escalabilidade
- Compreender as propriedades transacionais para os BDs NO-SQL
- Estabelecer as diferenças destas propriedades para BDs Relacionais e BDs NO-SQL
- Compreender o Teorema CAP
- Compreender as combinações CAP para os sistemas NO-SQL e estabelecer suas diferenças e aplicações



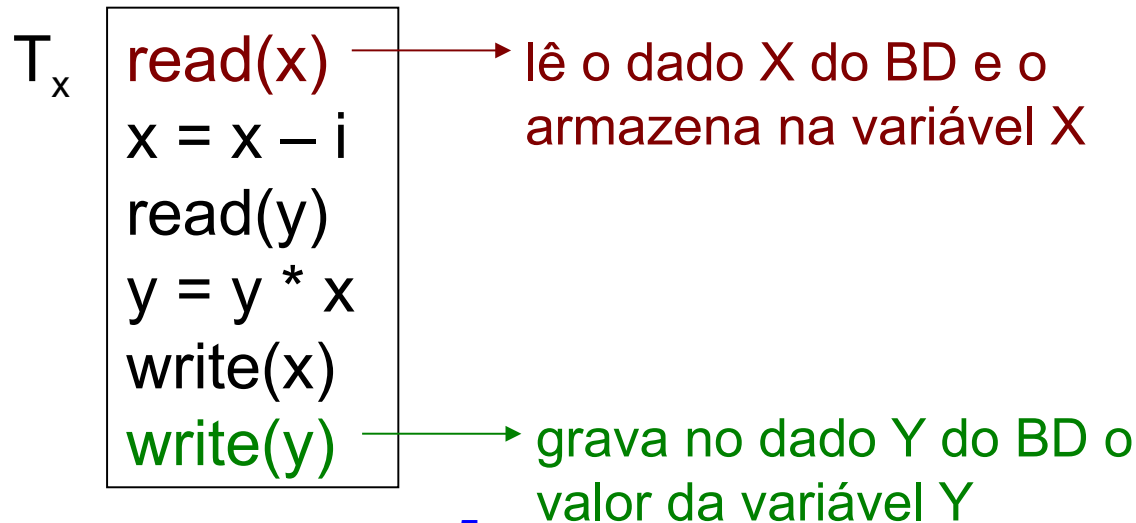
# 1- Introdução a Transações

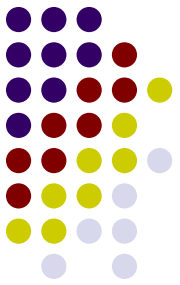
- SGBD
  - sistema de processamento de operações de acesso ao BD
- SGBDs são em geral **multi-usuários**
  - processam simultaneamente operações disparadas por vários usuários
    - deseja-se alta disponibilidade e tempo de resposta pequeno
  - execução intercalada de conjuntos de operações
    - exemplo: enquanto um processo  $i$  faz I/O, outro processo  $j$  é selecionado para execução
- Operações são chamadas **transações**



# 1.1 – Conceito de Transação

- Unidade lógica de processamento em um SGBD
- Composta de uma ou mais operações
  - seus limites podem ser determinados em SQL
- De forma abstrata e simplificada, uma transação pode ser encarada como um conjunto de operações de leitura e escrita de dados

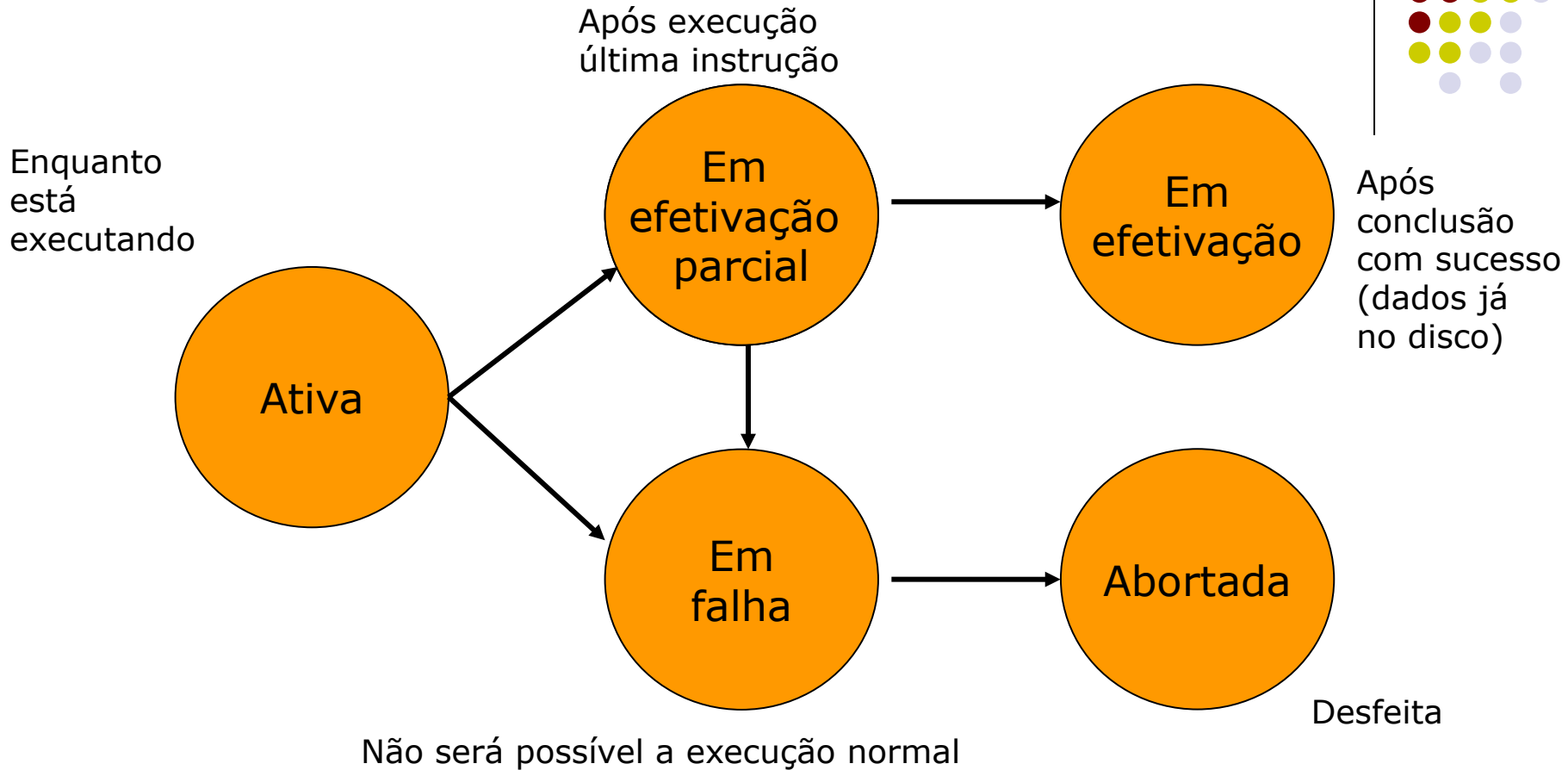
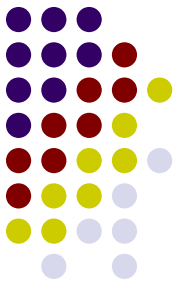




# 1.2 - Estados de uma Transação

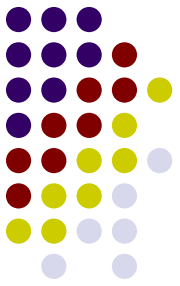
- Uma transação é sempre monitorada pelo SGBD quanto ao seu estado
  - que operações já fez? concluiu suas operações? deve abortar?
- Estados de uma transação
  - Ativa, Em processo de efetivação, Efetivada, Em falha, Em Efetivação(Concluída)
  - Respeita um Grafo de Transição de Estados

# 1.2.1 - Grafo de estados



- A transação está concluída se estiver em efetivação ou abortada
- Em efetivação parcial: ainda na memória
- Supõe-se, por enquanto, que falhas não resultam em perdas no disco

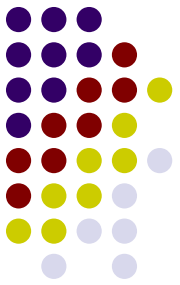
# 1.3 - Propriedades de uma Transação nos SGBDs Relacionais



- Requisitos que sempre devem ser atendidos por uma transação
- Chamadas de **propriedades ACID**
  - **A**tomicidade
  - **C**onsistência
  - **I**solamento
  - **D**urabilidade ou Persistência



# 1.4 - Atomicidade



- Princípio do ***“Tudo ou Nada”***
  - ou todas as operações da transação são efetivadas (**commit**) com sucesso no BD ou nenhuma delas se efetiva (**rollback**)
    - preservar a integridade do BD
- Responsabilidade do subsistema de recuperação contra falhas (**subsistema de recovery**) do SGBD
  - **desfazer (rollback)** as ações de transações parcialmente executadas

# 1.4 – Atomicidade (2)



- Deve ser garantida, pois uma transação pode manter o BD em um estado inconsistente durante a sua execução

Contas

número	saldo
100	500.00
200	200.00
...	

← x

← y

execução ↓

$T_x$  (transferência bancária)

read(x)

x.saldo = x.saldo – 100.00

write(x)

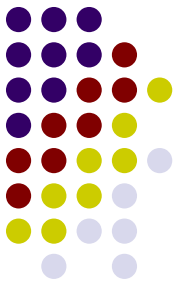
read(y)

y.saldo = y.saldo + 100.00

write(y)

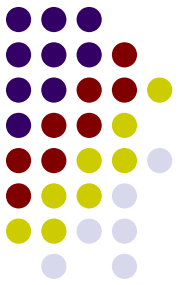
← falha!

# 1.5 - Consistência



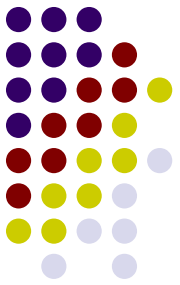
- Uma transação sempre conduz o BD de um estado consistente para outro estado também consistente
- Responsabilidade conjunta do
  - DBA
    - definir todas as **RIs** para garantir estados e transições de estado válidos para os dados
      - exemplos:  $\text{salário} > 0$ ;  $\text{salário novo} > \text{salário antigo}$
  - subsistema de *recovery*
    - desfazer as ações da transação que violou a integridade

# 1.6 - Isolamento



- No contexto de um conjunto de transações concorrentes, a execução de uma transação  $T_x$  deve funcionar como se  $T_x$  executasse de forma isolada
  - $T_x$  não deve sofrer interferências de outras transações executando concorrentemente
- Responsabilidade do subsistema de controle de concorrência (*scheduler*) do SGBD
  - garantir escalonamentos sem interferências

# 1.6 – Isolamento (2)



T <sub>1</sub>	T <sub>2</sub>
read(A) A = A – 50 write(A)	read(A) A = A+A*0.1 write(A)
read(B) B = B + 50 write(B)	read(B) B = B - A write(B)

escalonamento válido

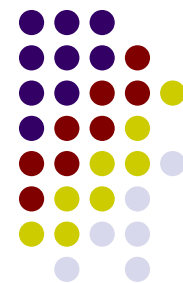
T <sub>1</sub>	T <sub>2</sub>
read(A) A = A – 50	read(A) A = A+A*0.1 write(A) read(B)
write(A) ←	
read(B) B = B + 50 write(B)	
	B = B - A write(B) ←

T<sub>1</sub> interfere em T<sub>2</sub>

T<sub>2</sub> interfere em T<sub>1</sub>

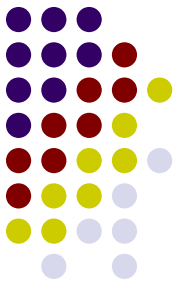
escalonamento inválido

# 1.7 - Durabilidade



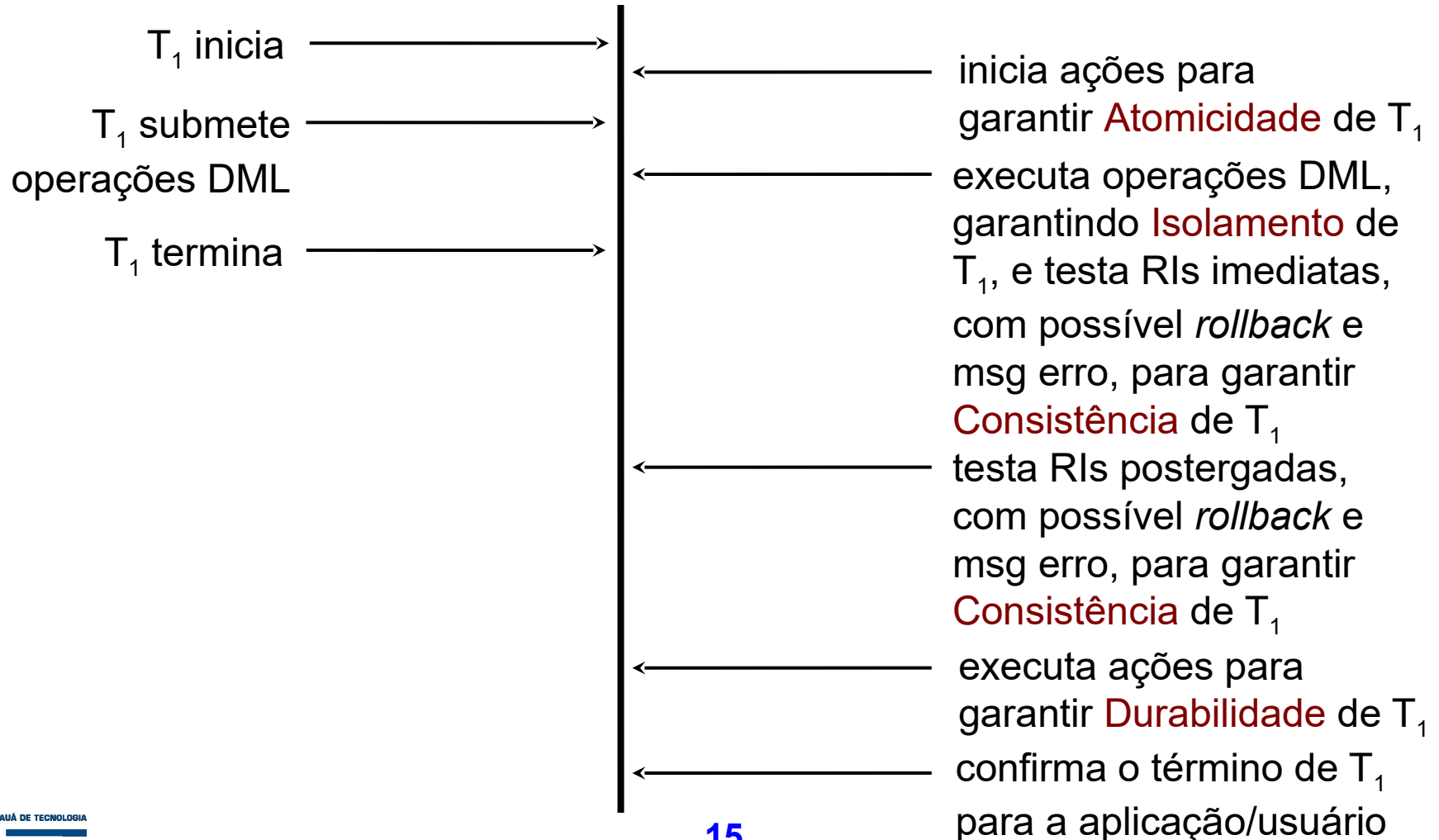
- Deve-se garantir que as **modificações realizadas por uma transação que concluiu com sucesso persistam no BD**
  - nenhuma falha posterior ocorrida no BD deve perder essas modificações
- Responsabilidade do **subsistema de *recovery***
  - **refazer** transações que executaram com sucesso em caso de falha no BD

# 1.8- Gerência Básica de Transações

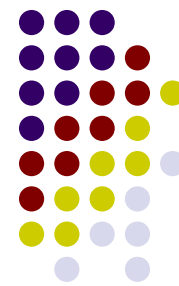


## Ações da Aplicação ou Usuário

## Ações do SGBD



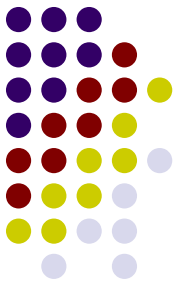
# 1.9 - Transações em SQL



- Por *default*, todo comando individual é considerado uma transação
  - exemplo: `DELETE FROM Pacientes`
    - exclui todas ou não exclui nenhuma tupla de pacientes, deve manter o BD consistente, etc
- SQL Padrão (SQL-92)
  - `SET TRANSACTION`
    - inicia e configura características de uma transação
  - `COMMIT [WORK]`
    - encerra a transação (solicita efetivação das suas ações)
  - `ROLLBACK [WORK]`
    - solicita que as ações da transação sejam desfeitas



# 1.9 - Transações em SQL (2)



- Principais configurações (**SET TRANSACTION**)
  - modo de acesso
    - **READ** (somente leitura), **WRITE** (somente atualização) ou **READ WRITE** (ambos - *default*)
  - nível de isolamento
    - indicado pela cláusula **ISOLATION LEVEL** *nível*
    - *nível* para uma transação  $T_i$  pode assumir
      - **SERIALIZABLE** ( $T_i$  executa com completo isolamento - *default*)
      - **REPEATABLE READ** ( $T_i$  só lê dados efetivados e outras transações não podem escrever em dados lidos por  $T_i$ ) – pode ocorrer que  $T_i$  só consiga ler alguns dados que deseja
      - **READ COMMITTED** ( $T_i$  só lê dados efetivados, mas outras transações podem escrever em dados lidos por  $T_i$ )
      - **READ UNCOMMITTED** ( $T_i$  pode ler dados que ainda não sofreram efetivação)

# 1.9 - Transações em SQL (3)

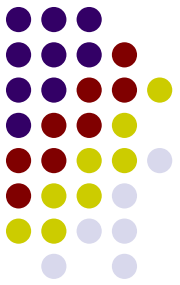
- Exemplo

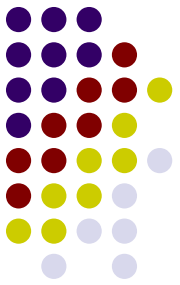
```
EXEC SQL SET TRANSACTION
        WRITE
        ISOLATION LEVEL SERIALIZABLE;

...
for (;;)
{
...
EXEC SQL INSERT INTO Empregados
        VALUES (:ID, :nome, :salario)

...
EXEC SQL UPDATE Empregados
        SET salário = salário + 100.00
        WHERE ID = :cod_emp
if (SQLCA.SQLCODE <= 0) EXEC SQL ROLLBACK;
...
}
EXEC SQL COMMIT;

...
```

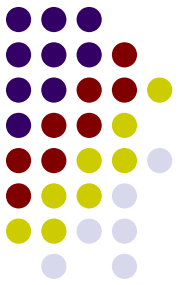




# 1.10 - Execução Concorrente

- Permitir que múltiplas transações concorram na atualização de dados traz diversas complicações em relação à consistência desses dados.
- Razões para permitir a concorrência:
  - Transação consiste em diversos passos :
    - atividades de I/O
    - atividades de CPU
    - Podem ser executadas em paralelo -> aumentar o rendimento
  - Mistura de transações em execução simultânea : algumas curtas e outras longas. Se a execução das transações for seqüencial, uma transação curta pode ser obrigada a esperar até que uma transação longa precedente se complete. Assim reduz-se o tempo médio de resposta: o tempo médio para uma transação ser completada após ser submetida.

# 1.10- Execução concorrente (2)



T1: transfere fundos de A  
para B

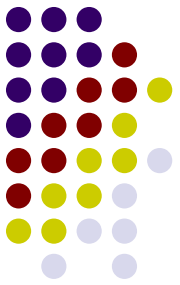
```
read(A);  
A := A - 50;  
write(A);  
read(B);  
B := B + 50;  
write(B);
```

T2: transfere 10% de A  
para B

```
read(A);  
temp := A * 0,1;  
A := A - temp;  
write(A);  
read(B);  
B := B + temp;  
write(B);
```

↓ tempo

# 1.10- Execução concorrente (3)



Escalas de execução (*schedule*) em seqüência: observe que o estado do BD é sempre consistente.

T1	T2	T1	T2
<pre>read(A); A := A - 50; write(A); read(B); B := B + 50; write(B);</pre>	<pre>read(A); temp := A * 0,1; A := A - temp; write(A); read(B); B := B + temp; write(B);</pre>	<pre>read(A); A := A - 50; write(A); read(B); B := B + 50; write(B);</pre>	<pre>read(A); temp := A * 0,1; A := A - temp; write(A); read(B); B := B + temp; write(B);</pre>

# 1.10 - Execução concorrente (4)



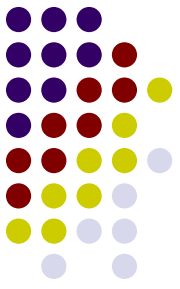
Correta

Incorreta

T1	T2
<code>read(A); A := A - 50; write(A);</code>	<code>read(A); temp := A * 0,1; A := A - temp; write(A);</code>
<code>read(B); B := B + 50; write(B);</code>	<code>read(B); B := B + temp; write(B);</code>

T1	T2
<code>read(A); A := A - 50;</code>	<code>read(A); temp := A * 0,1; A := A - temp; write(A);</code>
<code>write(A); read(B); B := B + 50; write(B);</code>	<code>read(B); B := B + temp; write(B);</code>

# 1.11 - Problema das atualizações perdidas

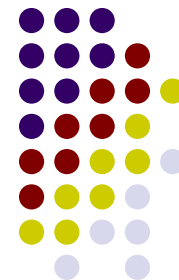


- Ocorre quando duas transações que acessam o mesmo item de dados possuem suas operações intercaladas de tal forma que o valor de algum item de dados fique incorreto.

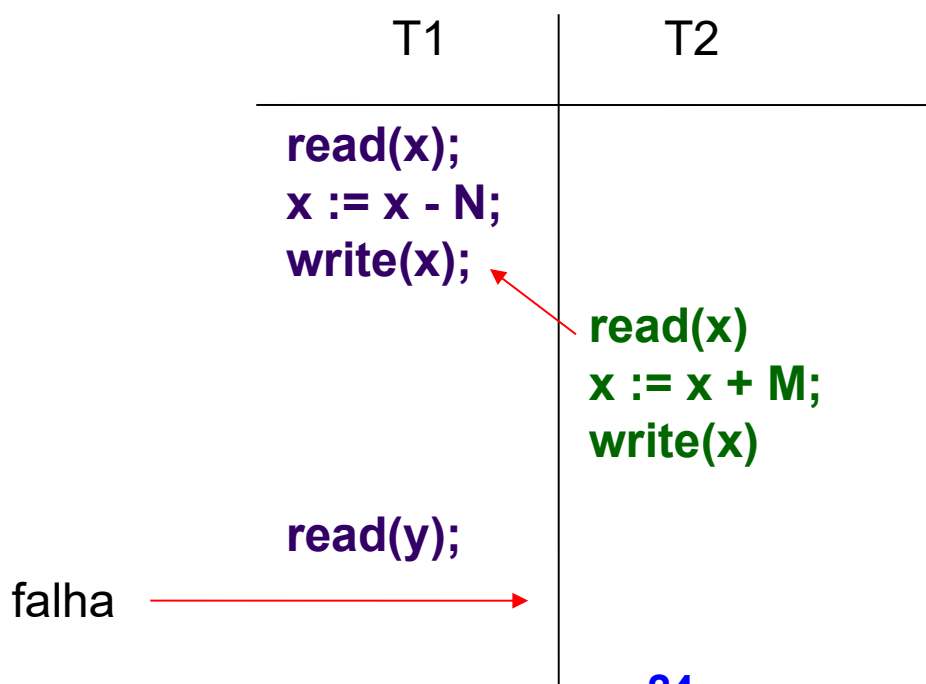
T1	T2
read(x); x := x - N;	read(x) x := x + M;
write(x); read(y);	write(x)
y := y + N; write(y)	



# 1.12 - Problema da dependência sem Commit (leitura suja)

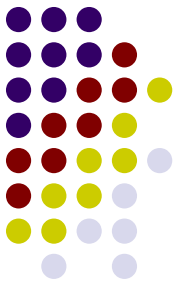


- Ocorre quando uma transação altera um item de dados e depois ela falha por alguma razão. O item de dado é acessado por outra transação antes que o valor original seja confirmado.

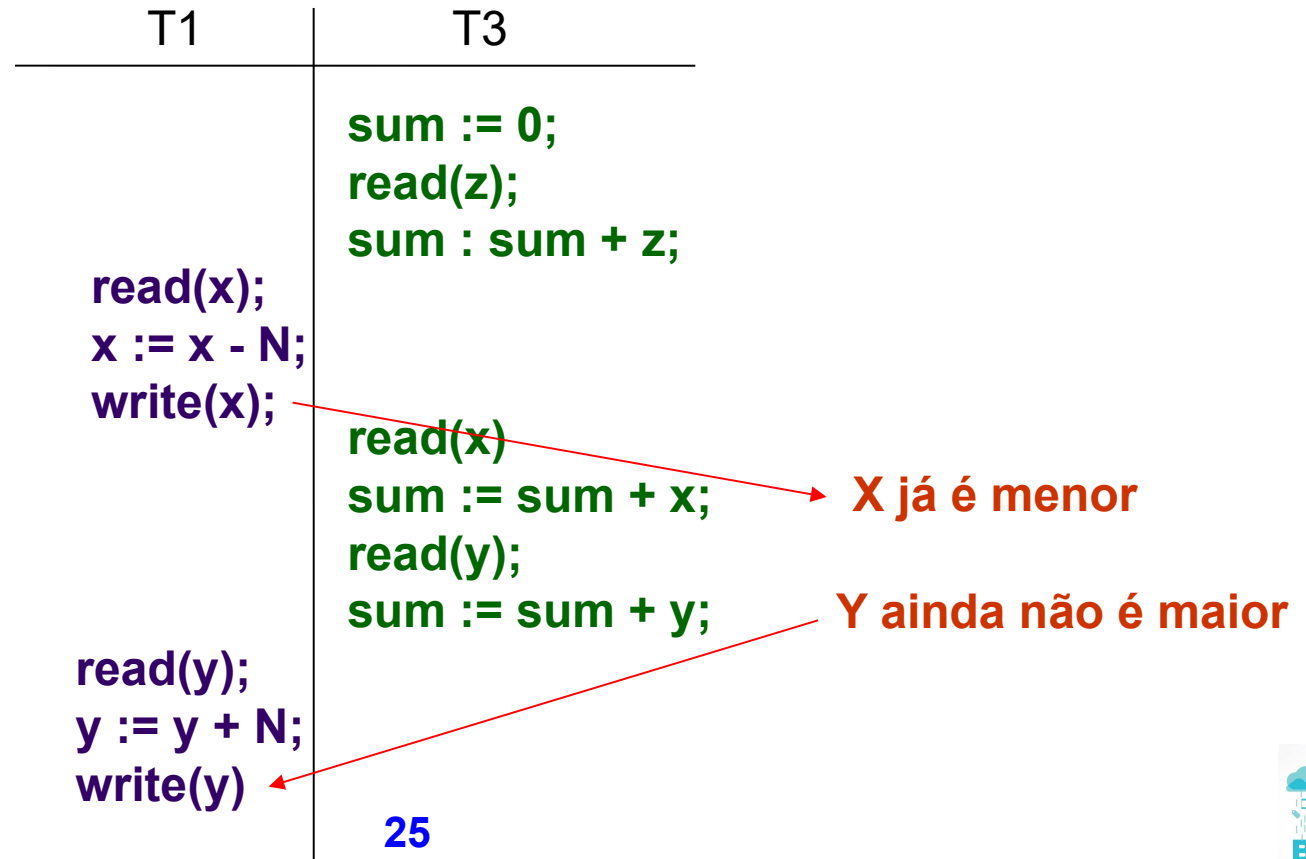




# 1.13 - Problema do resumo incorreto (análise inconsistente)



- Se uma transação está calculando uma função agregada com um conjunto de registros e outras transações estão alterando alguns destes registros a função agregada pode calcular alguns valores antes deles serem alterados e outros depois de serem alterados.

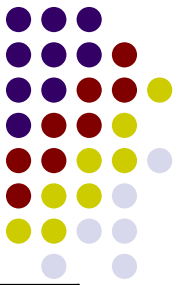


# 2 - Serialização



- O sistema gerenciador de banco de dados deve controlar a execução concorrente de transações para assegurar que o estado do banco de dados permaneça consistente.
- A consistência do banco de dados, sob execução concorrente, pode ser assegurada garantindo-se que qualquer escala executada concorrentemente tenha o mesmo efeito de outra que tivesse sido executada sem qualquer concorrência.
- Isto é, uma escala de execução deve, de alguma forma, ser **equivalente a uma escala seqüencial (transações sem intercalação, executadas de forma serial)** .
- Formas de equivalência entre escalas de execução podem ser verificadas sob duas visões:
  - Por conflito
  - Por visão

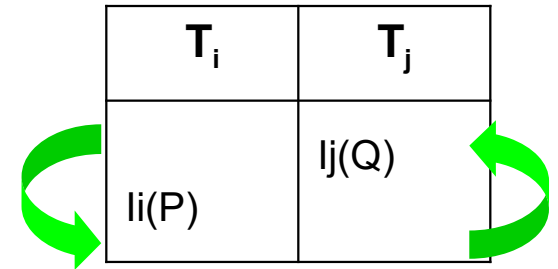
# 2.1 - Serialização por Conflito (1)



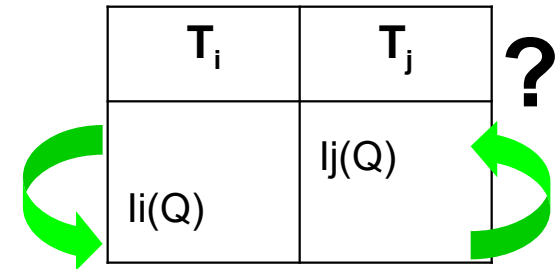
- Considere uma escala de execução  $S$  com duas instruções sucessivas,  $I_i$  e  $I_j$ , das transações  $T_i$  e  $T_j$  ( $i \neq j$ ), respectivamente.

$T_i$	$T_j$
$I_i(P)$	$I_j(Q)$

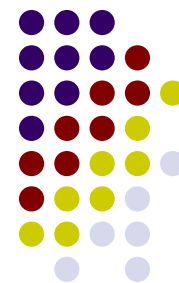
- Se  $I_i$  e  $I_j$  referem-se a itens de dados diferentes, então é permitido alternar  $I_i$  e  $I_j$  sem afetar os resultados de qualquer instrução da escala.



- Se  $I_i$  e  $I_j$  referem-se ao mesmo item de dados  $Q$ , então a ordem dos dois passos pode importar.



## 2.1 - Serialização por Conflito (2)



### 1- $I_i = \text{read}(Q)$ e $I_j = \text{read}(Q)$

A seqüência de execução de  $I_i$  e  $I_j$  não importa, já que o mesmo valor de  $Q$  é lido por  $T_i$  e  $T_j$ , independentemente da ordem destas operações.

### 2- $I_i = \text{read}(Q)$ e $I_j = \text{write}(Q)$

Se  $I_i$  vier antes de  $I_j$ , então  $T_i$  não lê o valor de  $Q$  que é escrito por  $T_j$  na instrução  $I_j$ .

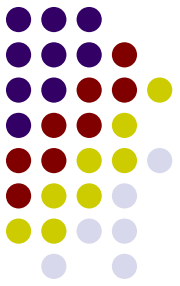
Se  $I_j$  vier antes de  $I_i$ , então  $T_i$  lê o valor de  $Q$  que é escrito por  $T_j$ .

Assim, a ordem de  $I_i$  e  $I_j$  importa.

### 3. $I_i = \text{write}(Q)$ e $I_j = \text{read}(Q)$

A ordem de  $I_i$  e  $I_j$  importa por razões semelhantes às do caso anterior.

## 2.1 - Serialização por Conflito (3)



4.  $I_i = \text{write}(Q)$  e  $I_j = \text{write}(Q)$

Como ambas as instruções são operações de escrita, a ordem dessas instruções não afeta  $T_i$  ou  $T_j$ .

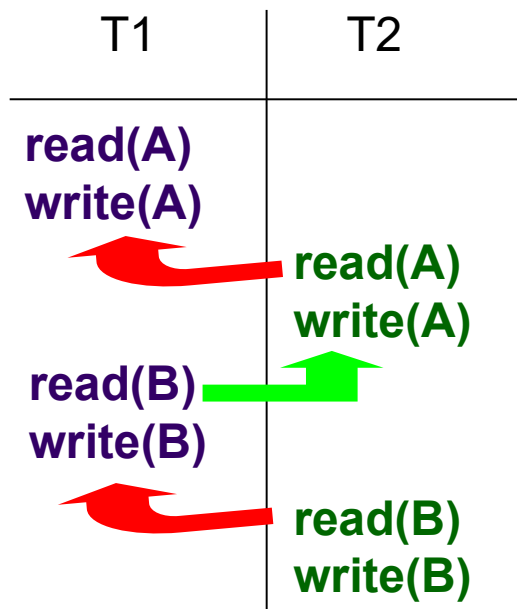
Entretanto, o valor obtido pela próxima instrução **read(Q)** em  $S$  é afetado, já que somente o resultado da última das duas instruções **write(Q)** é preservado no banco de dados.

Se não houver nenhuma outra instrução de **write(Q)** depois de  $I_i$  e  $I_j$  em  $S$ , então a ordem de  $I_i$  e  $I_j$  afeta diretamente o valor final de  $Q$  no que se refere ao estado do banco de dados após a execução da escala  $S$ .

## 2.1 - Serialização por Conflito (4)



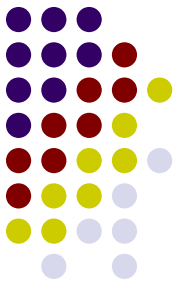
- Diz-se que duas instruções entram em **conflito** se elas são operações pertencentes a **transações diferentes, agindo no mesmo item de dados**, e pelo menos uma dessas instruções é uma operação de escrita (**write**).



A instrução `write(A)` de T1 entra em conflito com a instrução `read(A)` de T2.

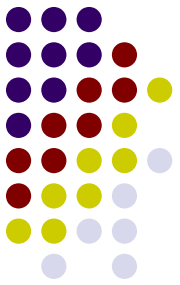
Porém, a instrução `write(A)` de T2 não está em conflito com a instrução `read(B)` de T1.

## 2.2 - Serialização por Conflito (escalas equivalentes)



- Sejam  $I_i$  e  $I_j$  instruções consecutivas de uma escala de execução  $S$ .
- Se  $I_i$  e  $I_j$  são instruções de transações diferentes e não entram em conflito, então podemos trocar a ordem de  $I_i$  e  $I_j$  para produzir uma nova escala de execução  $S'$ .
- Diz-se que  $S$  e  $S'$  são equivalentes já que todas as instruções aparecem na mesma ordem em ambas as escalas de execução com exceção de  $I_i$  e  $I_j$ , cuja ordem não importa.

## 2.2 - Serialização por Conflito (escalas equivalentes (2))



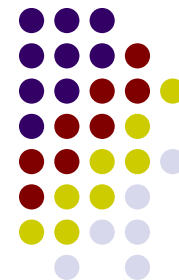
- Voltando à escala S anterior, a instrução `write(A)` de T2 não entra em conflito com a instrução `read(B)` de T1. Então é permitido trocar essas instruções para gerar uma escala de execução equivalente.

T1	T2
<code>read(A)</code> <code>write(A)</code>	
<code>read(B)</code> ↑	<code>read(A)</code> <code>write(A)</code> ↓
<code>write(B)</code>	<code>read(B)</code> <code>write(B)</code>

Em relação a um mesmo estado inicial do sistema, ambas as escalas (S e esta) Produzem o mesmo estado final no sistema.



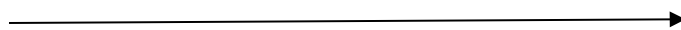
## 2.2 - Serialização por Conflito (escalas equivalentes (3))



- Se as seguintes trocas de instruções não-conflitantes forem feitas ...

- read(B) de T1 por read(A) de T2;
- write(B) de T1 por write(A) de T2;
- Write(B) de T1 por read(A) de T2

- ... uma escala com execuções seriais será obtida.



- Assim mostrou-se que a escala S é equivalente, no conflito, a uma escala seqüencial.
- 
- Essa equivalência garante que para um mesmo estado inicial do sistema, a escala S produzirá o mesmo estado final produzido por essa escala seqüencial.

**S'**

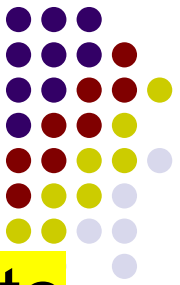
T1	T2
read(A) write(A) read(B) write(B)	read(A) write(A) read(B) write(B)

## 2.2 - Serialização por Conflito (escalas equivalentes (4))

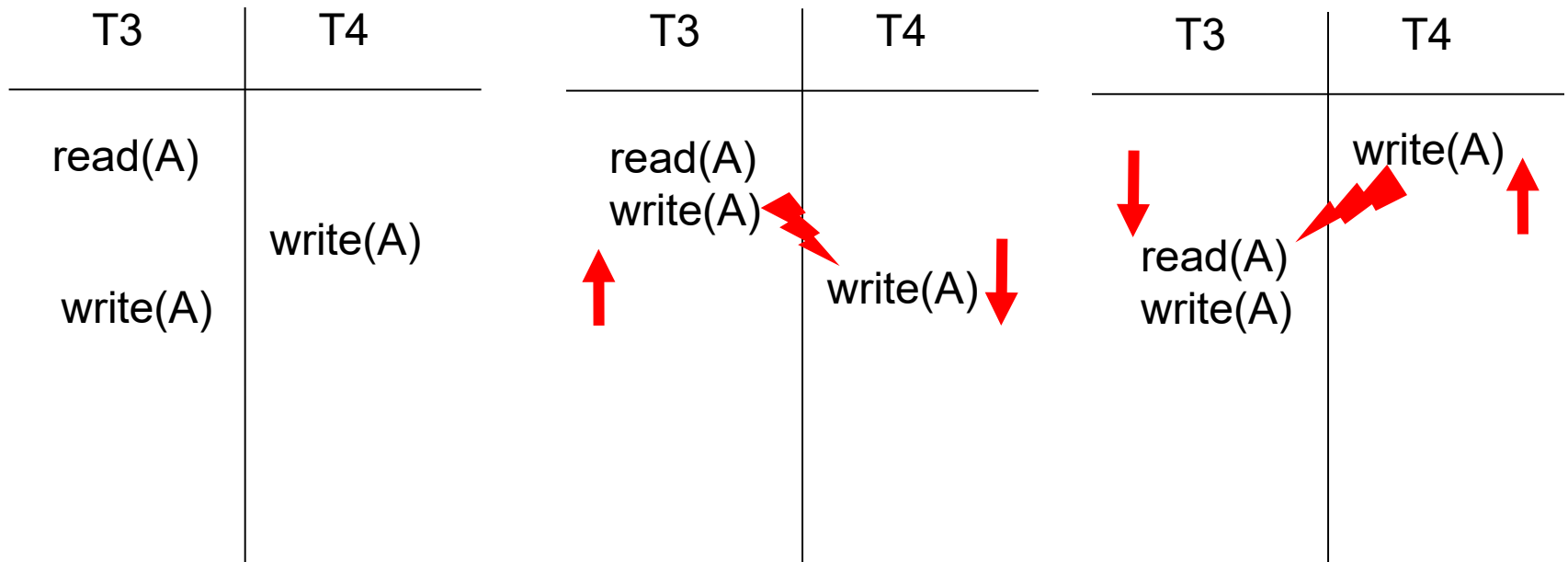


- Se uma escala de execução  $S$  puder ser transformada em outra,  $S'$ , por uma série de trocas de instruções não conflitantes, dizemos que  $S$  e  $S'$  são **equivalentes em conflito**.
- O conceito de equivalência em conflito leva ao conceito de serialização por conflito.
- Uma escala de execução  $S$  é serializável por conflito se ela é equivalente em conflito a uma escala de execução seqüencial.

## 2.2.1 - Serialização por conflito - Exemplo

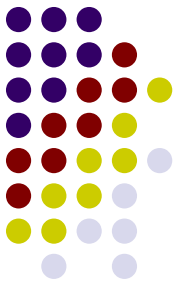
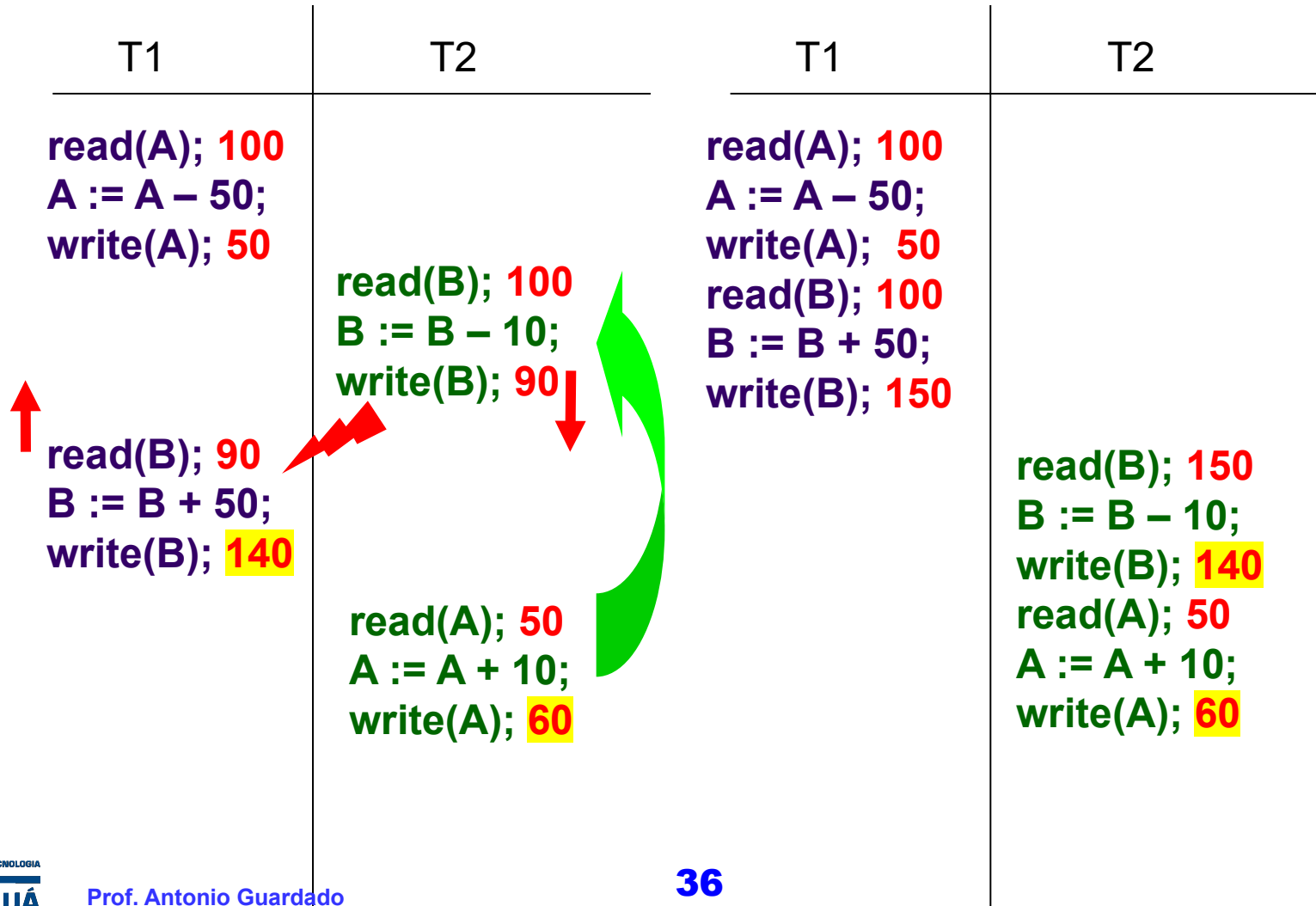


- A escala abaixo não é serializável por conflito pois não é equivalente em conflito nem à escala seqüencial  $\langle T3, T4 \rangle$  nem  $\langle T4, T3 \rangle$ .

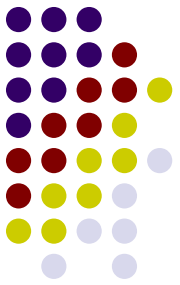


## 2.2.1 - Serialização por conflito – Exemplo

- A escala abaixo não é serializável por conflito, mas produz o mesmo resultado que uma escala seqüencial. Confira.



## 2.3 - Teste de serialização por conflito

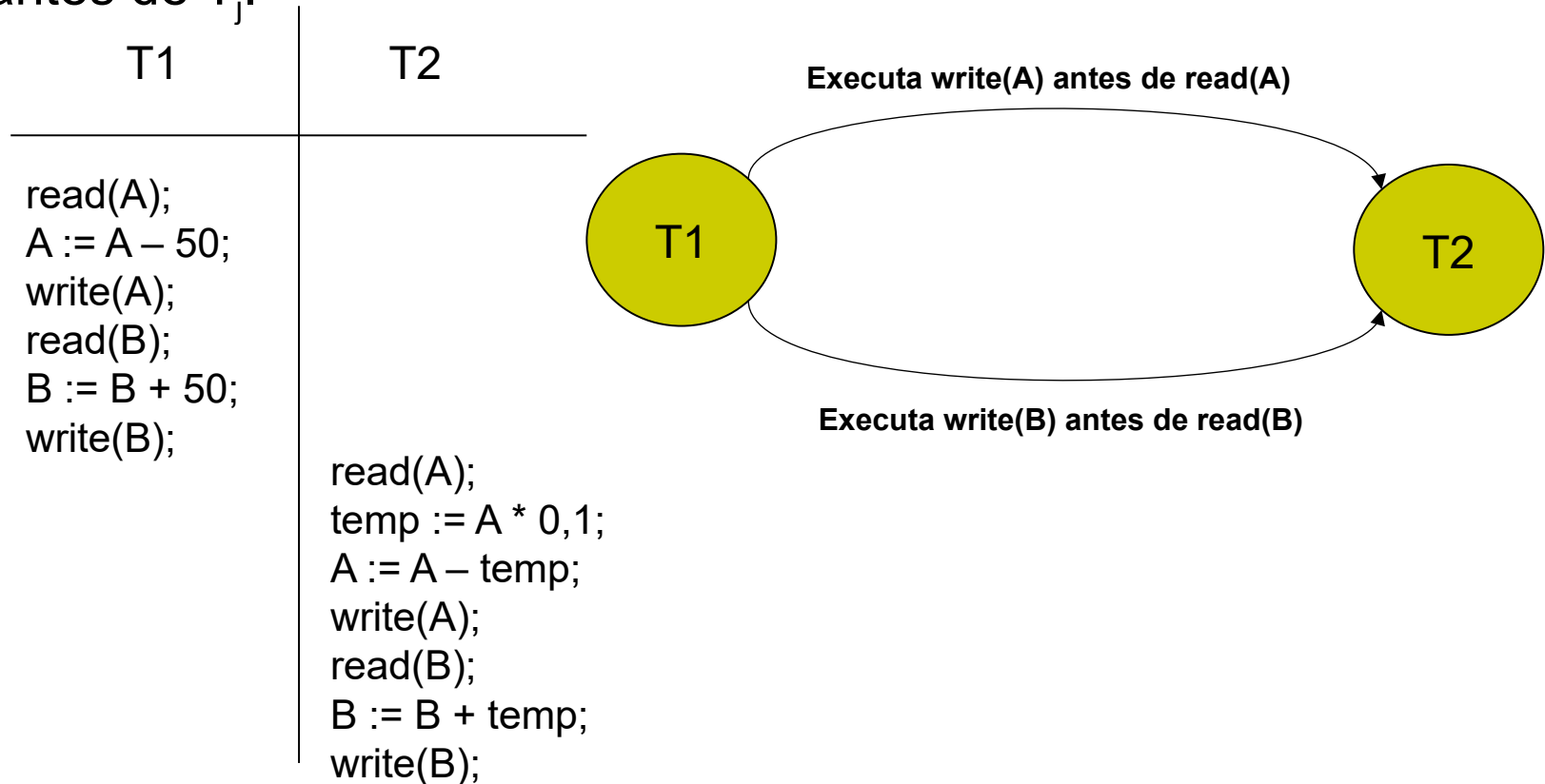


- Seja **S** uma escala. Para saber se ela é serializável em relação às operações conflitantes é necessário criar um **grafo de precedência** para S.
- $G = (V, E)$  em que  $V$  é um conjunto de vértices e  $E$  é um conjunto de arestas.
  - O conjunto de vértices é composto por todas as transações que participam da escala.
  - O conjunto de arestas consiste em todas as arestas  $T_i \rightarrow T_j$  para as quais uma das seguintes condições é verdadeira:
    - $T_i$  executa **write(Q)** antes de  $T_j$  executar **read(Q)**;
    - $T_i$  executa **read(Q)** antes de  $T_j$  executar **write(Q)**;
    - $T_i$  executa **write(Q)** antes de  $T_j$  executar **write(Q)**;

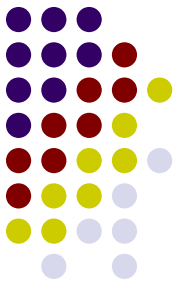
## 2.3 - Teste de serialização por conflito – Exemplo 1



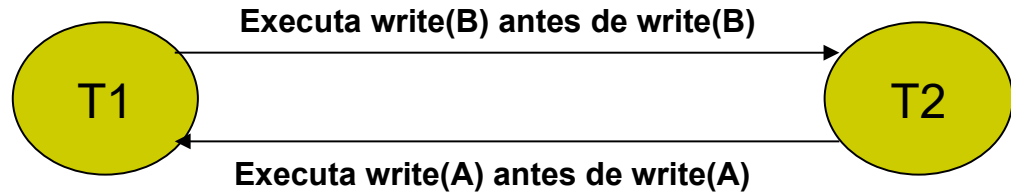
- Se existir uma aresta  $T_i \rightarrow T_j$  no grafo de precedência, então, em qualquer escala seqüencial  $S'$  equivalente a  $S$ ,  $T_i$  deve aparecer antes de  $T_j$ .



## 2.3 - Teste de serialização por conflito – Exemplo 2



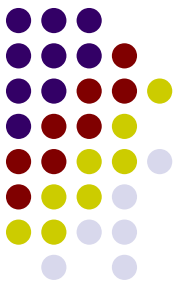
T1	T2
read(A) A := A - 50	read(A) temp := A * 0,1 A := A - temp write(A) read(B)
write(A) read(B) B := B + 50 write(B)	B := B + temp; write(B)



**Ciclo no grafo: se o grafo de precedência possui ciclo, então a escala S não é serializável por conflito.**

A ordem de serialização pode ser obtida por meio da classificação topológica, que estabelece uma ordem linear para a escala consistente com a ordem parcial do grafo de precedência.

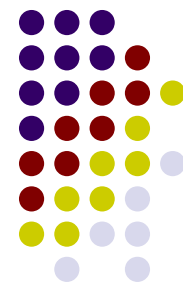
# 3 – Propriedades Transacionais em BDs NO-SQL



- Bancos de dados relacionais tem problemas
  - **Escalabilidade vertical**: adicionar recursos de hardware a uma mesma máquina ( + memória, + armazenamento, + processadores, + velocidade de processamento) ➡ limite físico computacional e alto custo
  - **Escalabilidade horizontal**: particionar os dados em várias máquinas, **torna a manutenção das tabelas bem difícil e complicado**; depende da fragmentação se é horizontal (linhas) ou vertical (colunas)



# 3 – Propriedades Transacionais em BDs NO-SQL



- Bancos de dados relacionais tem problemas
  - **alto volume de dados com baixa velocidade de resposta** - a medida que os dados crescem na base, fica cada vez mais difícil a disponibilização rápida das informações.
  - **os dados mudam constantemente** - alterar a entidade no mundo relacional em algumas situações não é fácil, dependendo do tipo do dado alguns bancos simplesmente **bloqueiam** o acesso a tabela para realizar uma simples alteração de coluna.

# 3 – Propriedades Transacionais em BDs NO-SQL



- **Bancos de dados relacionais tem problemas**

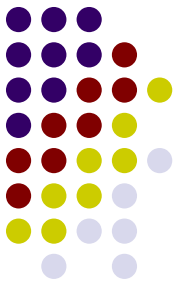
- Protocolos de bloqueio de efetivação em duas fases distribuído – 2PC – para preservar as propriedades ACID, em especial a consistência, aumentam consideravelmente o tempo de resposta do sistema



latência muito alta implica em **baixa disponibilidade**

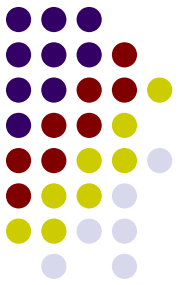
- Bancos de dados NO-SQL precisam de **alta disponibilidade e redundância** (mesmo dado em muitos nós) ➡ Não é possível usar propriedades ACID

# 4- Propriedades BASE - NoSQL

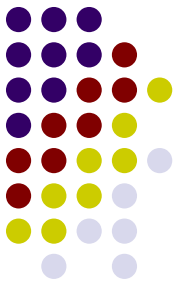


- **B**asically **A**vailable – **B**asicamente **D**isponível : **disponibilidade é prioridade**, o sistema deve estar em funcionamento na maior parte do tempo
- **S**oft-State – **E**stado **L**eve : **não precisa ser consistente o tempo todo**
- **E**ventually Consistent – **E**ventualmente **C**onsistente : **consistente em algum momento não determinado**, a consistência nem sempre é mantida para todos os nós, nós podem não ter a mesma versão dos dados.

# 4.1- Comparativo ACID x BASE



ACID	BASE
Consistência forte	Fraca consistência
Isolamento	Disponibilidade em primeiro lugar
Concentra-se em "commit"	Melhor esforço em disponibilidade para partições
Transações aninhadas	Respostas aproximadas
Conservador (pessimista, bloqueia todos os registros para evitar conflitos)	Agressivo (otimista, detectam os conflitos e depois faz o tratamento)
Evolução difícil (por exemplo, esquema)	Evolução mais fácil



# 5 – Teorema CAP

- **Consistência** – **C**onsistency.
- **Disponibilidade** – **A**vailability.
- **Tolerância ao Particionamento** - **P**artition tolerance.
- Teorema CAP : é impossível garantir essas **três propriedades ao mesmo tempo**
- é possível garantir quaisquer **duas dessas propriedades ao mesmo tempo**

Gilbert, S.; Lynch, N. A. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. SIGACT News, 33(2):51–59, 2002.

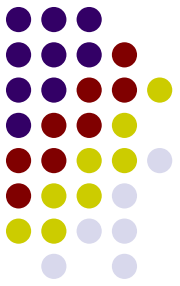
# 5.1 - Consistency - Consistência



O sistema **garante a leitura do dado mais atualizado**, quando ele foi escrito.

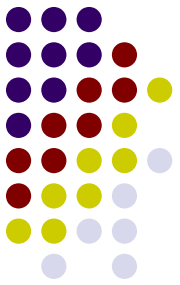
O cliente pode ler o dado no mesmo nó que este dado foi escrito ou de um nó diferente, o mesmo dado será retornado para a aplicação cliente. Mesmo que alguém tenha mudado o estado do domínio da aplicação com novas informações, o comportamento de consistência, garantirá que o cliente não verá dados velhos, apenas os novos dados serão visualizados.

## 5.2- Availability - Disponibilidade



- Quando o cliente lê ou escreve um dado em um dos nós, o nó que está sendo utilizado pelo cliente, pode estar indisponível.
- Isso não quer dizer que um nó não pode falhar, o que esse comportamento quer dizer é que, um nó pode sim ficar offline/cair/falhar, mas o **sistema/aplicação continuaria disponível**

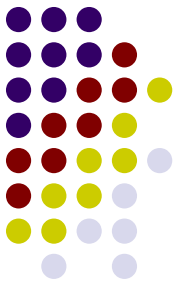
## 5.2- Availability - Disponibilidade



- Se um sistema é capaz de obter acesso de leitura/escrita em um nó que não possui falhas e este responde em um tempo razoável, temos aqui a garantia de disponibilidade.
- **Mas pode trazer um dado não atual**

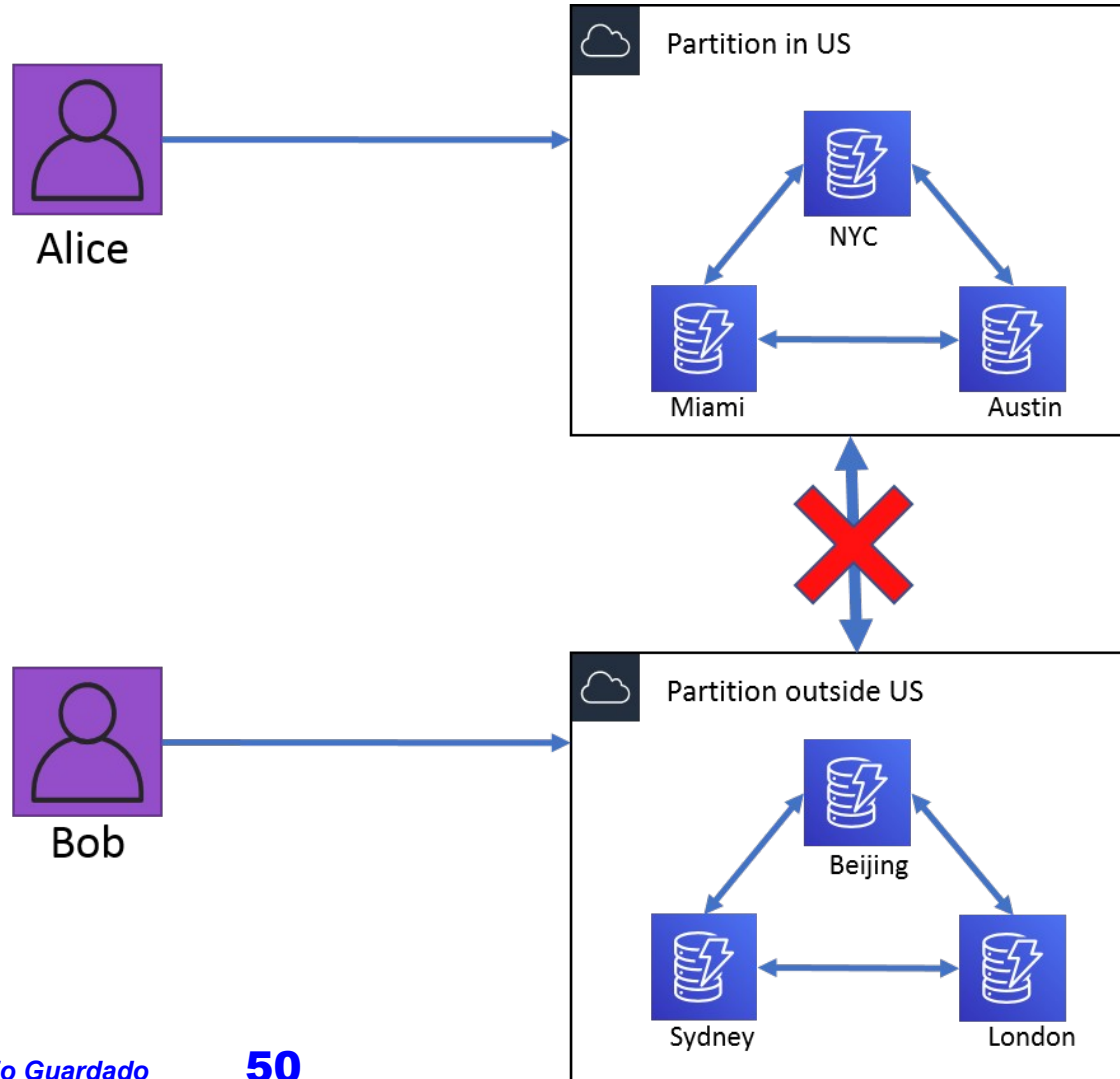
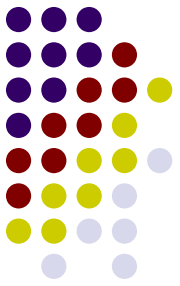


## 5.3- Partition Tolerance – Tolerância a falhas de partição

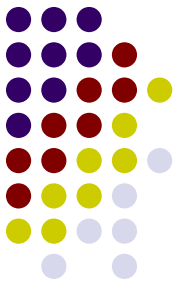


É a garantia de que o sistema continue operante mesmo no caso da ocorrência de uma falha que isole os nós em grupos, onde os nós de um grupo não consigam se comunicar com os dos demais grupos. O sistema **continua operando, mesmo que aconteça alguma falha na rede** (falha de conectividade).

# 5.3- Partition Tolerance – Tolerância a falhas de partição



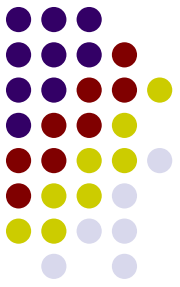
## 5.3- Partition Tolerance – Tolerância a falhas de partição



Os sistemas distribuídos são, por natureza, **não confiáveis**. Eles são suscetíveis a diferentes tipos de problemas: **falhas na rede, perda de mensagens, quebra de máquinas, ataques maliciosos, etc. Quanto mais nós no sistema, maior o risco de problemas.**

A ocorrência de problemas como esses podem impedir (ou dificultar) a comunicação entre os nós, causando uma **partição do sistema**: ou seja, **os nós podem ficar "isolados" em grupos que conseguem se comunicar internamente, mas que não conseguem se comunicar com outros grupos.**

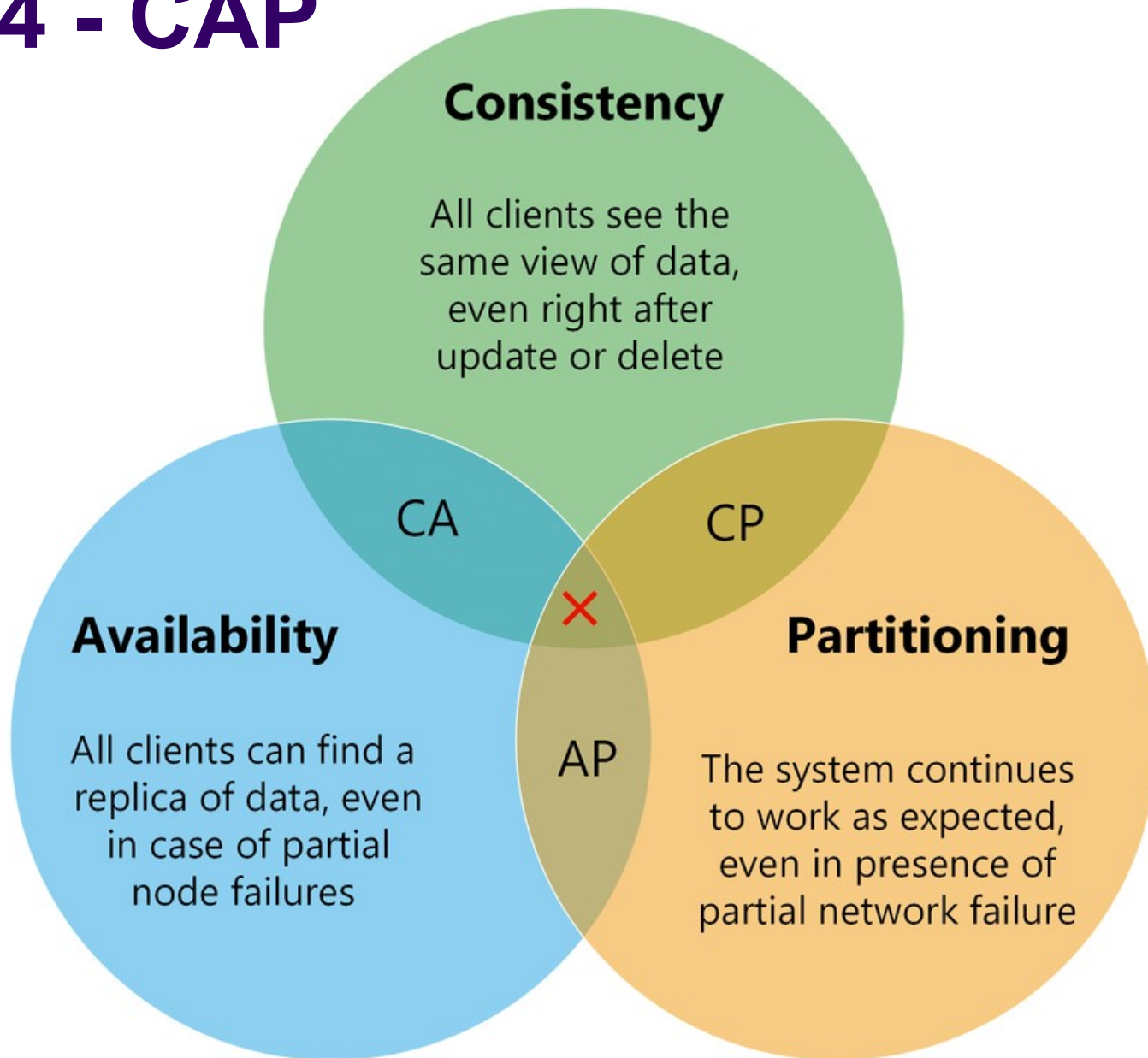
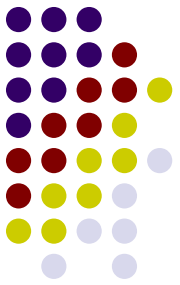
## 5.3- Partition Tolerance – Tolerância a falhas de partição



Para **garantir a consistência** em caso de partição do sistema, duas estratégias são possíveis:

- 1) fazer com que seus **nós deixem de receber requisições dos clientes** enquanto o problema da partição persistir;
- 2) **continuar recebendo as requisições dos clientes**, mas elas **só serão atendidas quando o problema da partição for resolvido.**

# 5.4 - CAP



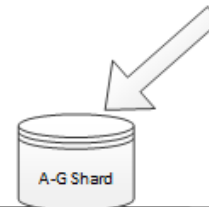
# 6 – Combinações CAP



- Dependendo da aplicação escolher entre consistência forte, alta disponibilidade e tolerância ao particionamento
- **Sistemas CA (Consistência e Disponibilidade)**
- **Sistemas CP (Consistência e Particionamento)**
- **Sistemas AP (Disponibilidade e Particionamento)**

# 6.1 – Modelos de Distribuição

Key	Name	Description	Stock	Price	LastOrdered
ARC1	Arc welder	250 Amps	8	119.00	25-Nov-2013
BRK8	Bracket	250mm	46	5.66	18-Nov-2013
BRK9	Bracket	400mm	82	6.98	1-Jul-2013
HOS8	Hose	1/2"	27	27.50	18-Aug-2013
WGT4	Widget	Green	16	13.99	3-Feb-2013
WGT6	Widget	Purple	76	13.99	31-Mar-2013

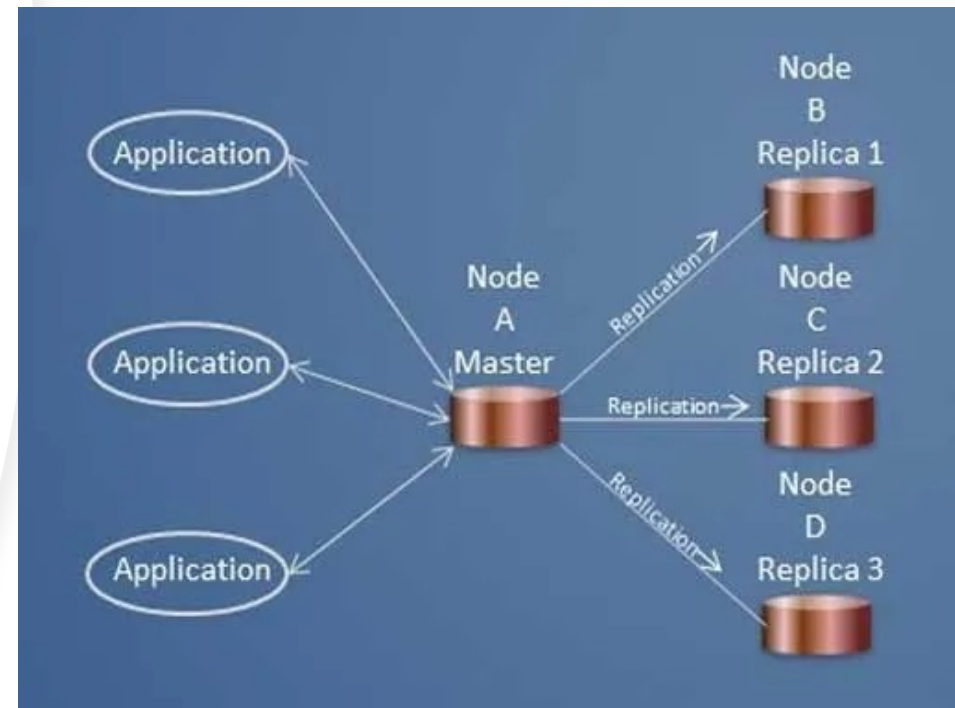


Key	Name	Description	Stock	Price	LastOrdered
ARC1	Arc welder	250 Amps	8	119.00	25-Nov-2013
BRK8	Bracket	250mm	46	5.66	18-Nov-2013
BRK9	Bracket	400mm	82	6.98	1-Jul-2013

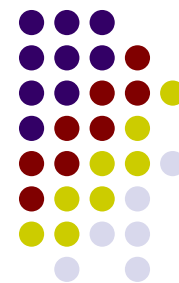


Key	Name	Description	Stock	Price	LastOrdered
HOS8	Hose	1/2"	27	27.50	18-Aug-2013
WGT4	Widget	Green	16	13.99	3-Feb-2013
WGT6	Widget	Purple	76	13.99	31-Mar-2013

- Particionamento (*Sharding*)
  - *colocar diferentes partes dos dados em diferentes máquinas*
  - *Permite paralelismo nos acessos*
  - Escala escritas e leituras
- Replicação
  - *Copiar os mesmos dados em diferentes locais*
  - Escala somente leituras
  - Dois modelos possíveis:
    - Mestre-Escravo e Par-a-Par
- Particionamento + Replicação



## 6.2 - Sistemas CA



Os sistemas com **consistência forte** e **alta disponibilidade** **não sabem lidar com a possível falha de uma partição.**

Caso ocorra, sistema inteiro pode ficar indisponível até o membro do cluster voltar.

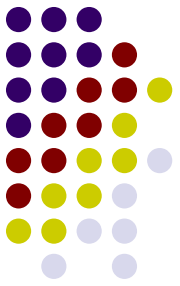
Situação em que o sistema garante **apenas alta disponibilidade (em um único nó).**

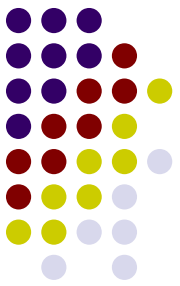
Uma máquina não pode ser particionada.  
-> BDs Relacionais Centralizados





## 6.2 - Sistemas CA





## 6.3 - Sistemas CP

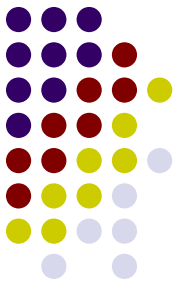
Para sistemas que precisam da **consistência forte e tolerância a particionamento** é necessário abrir a mão da disponibilidade (um pouco).

O sistema fica operante no caso de particionamento, mas pode demorar bastante tempo para conseguir responder às requisições que recebe ou, até mesmo, nunca respondê-las.

Exemplos são Google BigTable, HBase ou MongoDB entre vários outros.

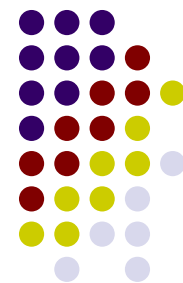


## 6.3.1- Consistência por meio de Quóruns



- Na replicação de dados, quanto mais nós envolvidos no tratamento de leituras ou escritas, maior a chance de evitar inconsistências
- **Quórum** de uma operação = quantidade de votos (confirmações) que a operação precisa receber para poder ser aplicada sobre um item de dados em um nó
- Cada nó com uma réplica do item tem direito a um voto

## 6.3.1- Quórum de Escrita

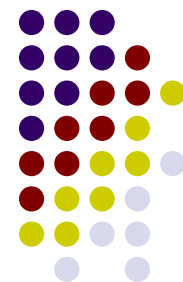


- $W$  = quórum de escrita
  - número de nós que confirmam uma escrita
- $N$  = fator de replicação dos dados
  - número de réplicas
- Para garantir escritas com consistência forte, a seguinte fórmula precisa ser respeitada:

$$W > N / 2$$

- O número de nós que confirmam uma escrita precisa ser maior que a metade do número de réplicas do item a ser escrito, ou seja, a maioria
- Dessa forma, se houver escritas conflitantes, somente uma será confirmada pela maioria dos nós

## 6.3.1- Quórum de Leitura

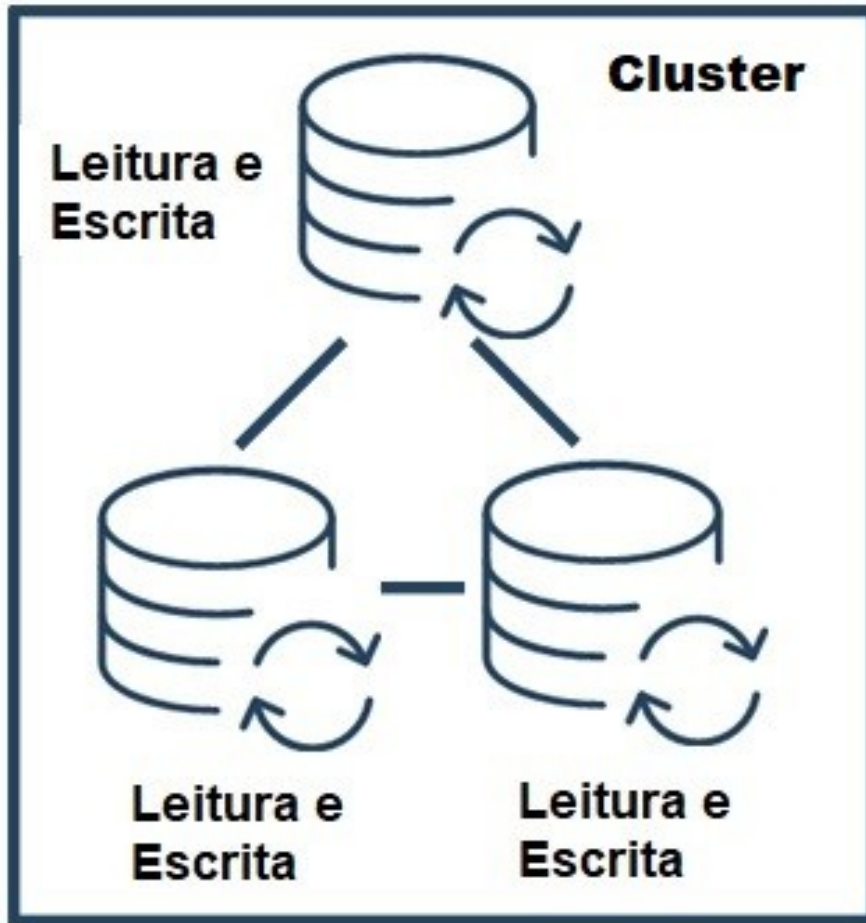
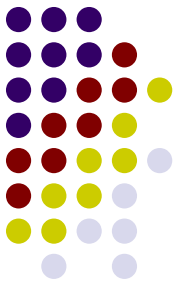


- Quantos nós é preciso contatar para se ter a garantia de que o valor mais recente para um item de dado foi lido?
  - A resposta depende do quórum de escrita  $W$
- $R$  = quórum de leitura
  - número de nós contatados para a leitura
- Para garantir leituras altamente consistentes, a seguinte fórmula precisa ser respeitada:

$$R + W > N$$

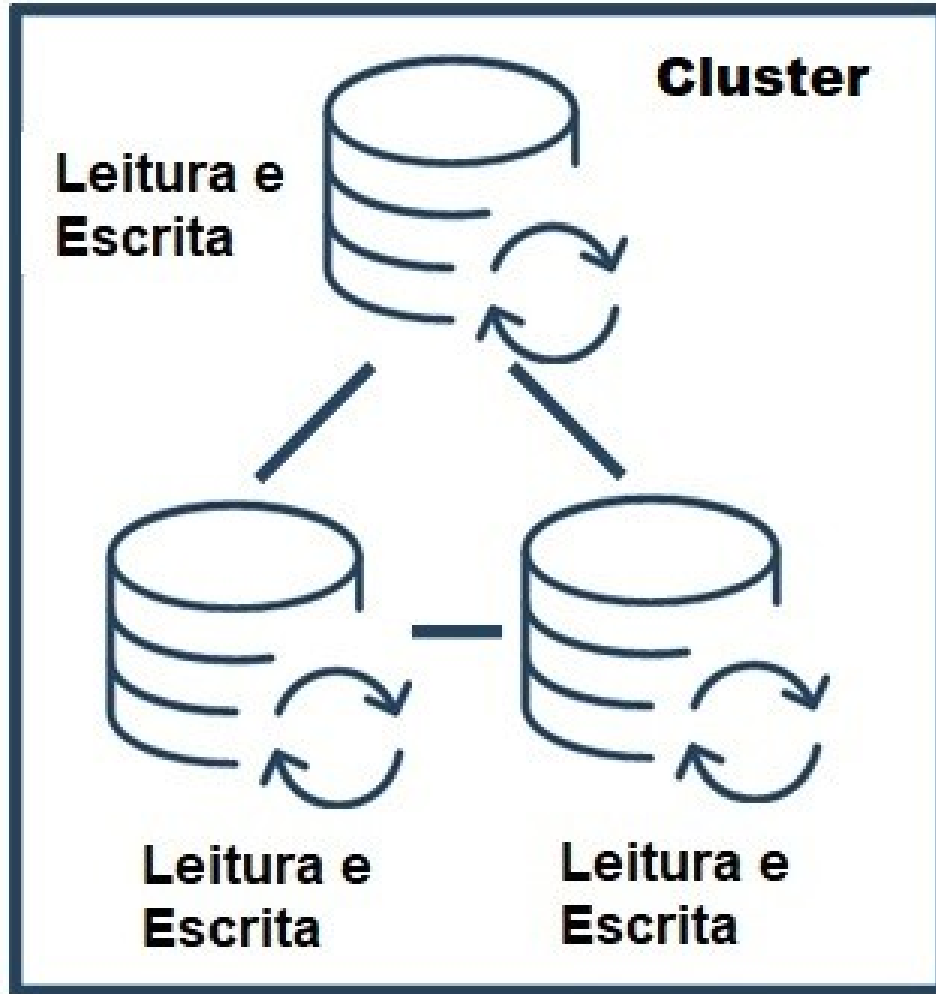
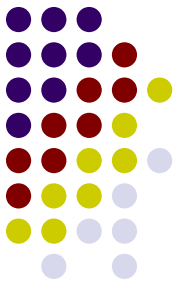
- A fórmula garante que, entre os  $R$  nós de réplicas contatados, haverá pelos menos 1 que tem a versão mais nova do item de dado

## 6.3.2 - Sistemas CP - Consenso



Se precisar de leituras rápidas e fortemente consistentes, pode-se exigir que as escritas sejam reconhecidas por todos os nós, permitindo assim que as leituras entrem em contato com apenas um. Isso significaria que as escritas são lentas, pois elas precisam entrar em contato com todos os nós, e o sistema não seria capaz de tolerar a perda de um nó.

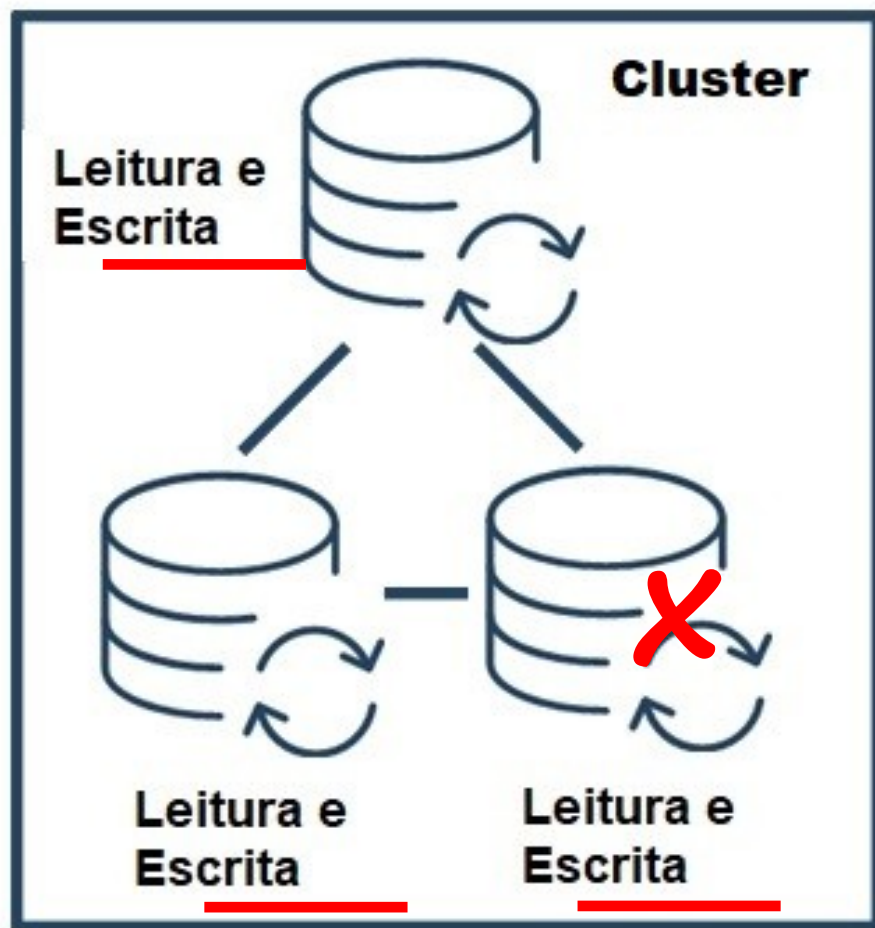
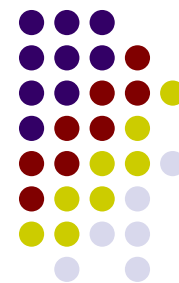
## 6.3.2 - Sistemas CP - Consenso



**Consenso** : Fator de Replicação  $N$  igual a 3  
Consideremos  $R = 1$  e  $W = 3$

Portanto, para escrever no banco é necessária a confirmação de escrita de três nós (quórum)  
 $W > N/2$  e  $R+W > N$

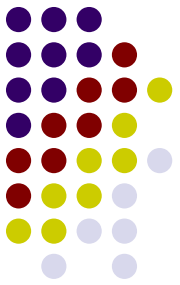
## 6.3.2 - Sistemas CP - Consenso



**Consenso** : Com a **falha de um nó a escrita é comprometida** (não há quórum, precisa  $W = 3$ )  
**Perde-se a disponibilidade (A) para escrita (sistema para de escrever)**, mas a **leitura** ainda pode ocorrer em qualquer um dos demais nós ( $R = 1$ ). Se **falhar mais um nó** a **leitura** ainda pode ser realizada com consistência.



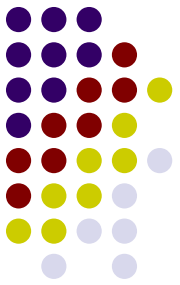
## 6.3.2 - Sistemas CP



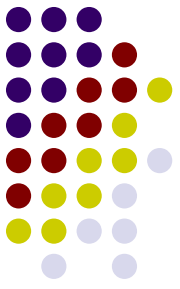
## 6.4 - Sistemas AP

• Há sistemas que jamais podem ficar offline, portanto não desejam sacrificar a disponibilidade. Para ter alta disponibilidade mesmo com uma tolerância a particionamento é preciso comprometer a consistência

• Exemplos de Bancos são: Cassandra, MongoDB, Voldemort

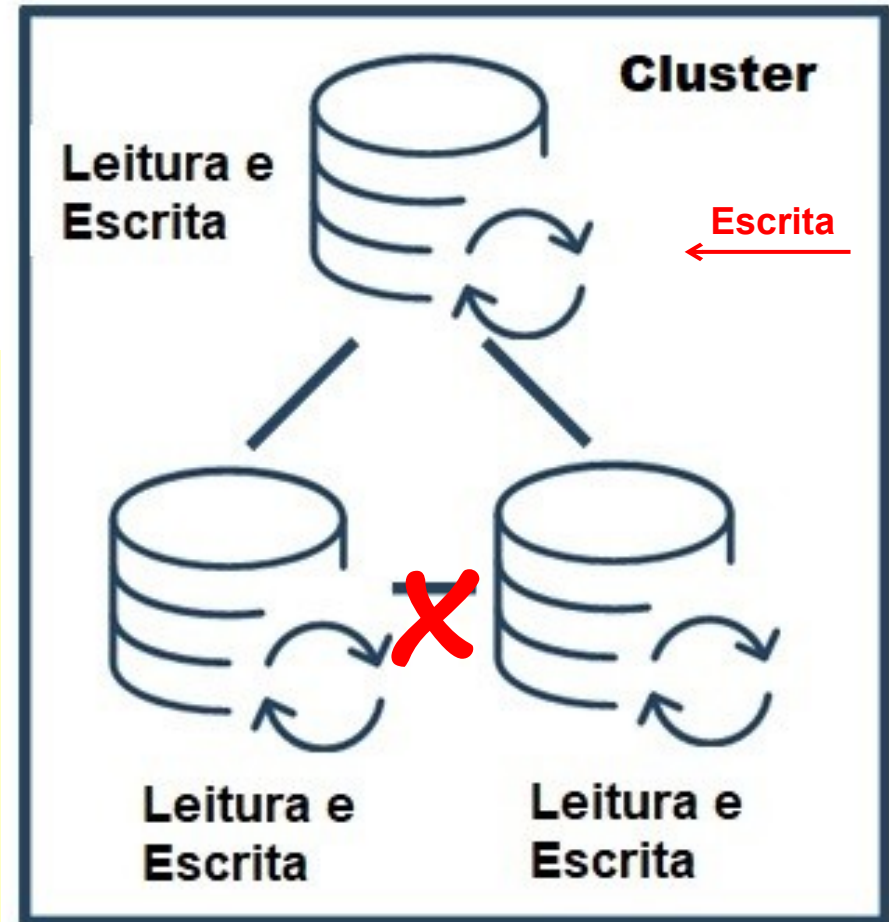


## 6.4 - Sistemas AP

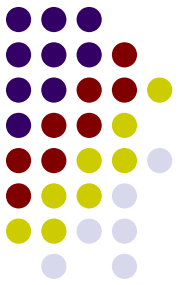


**Consistência Eventual :**  
sempre disponível para  
escrita e depois sincroniza  
os dados com os demais  
nós.

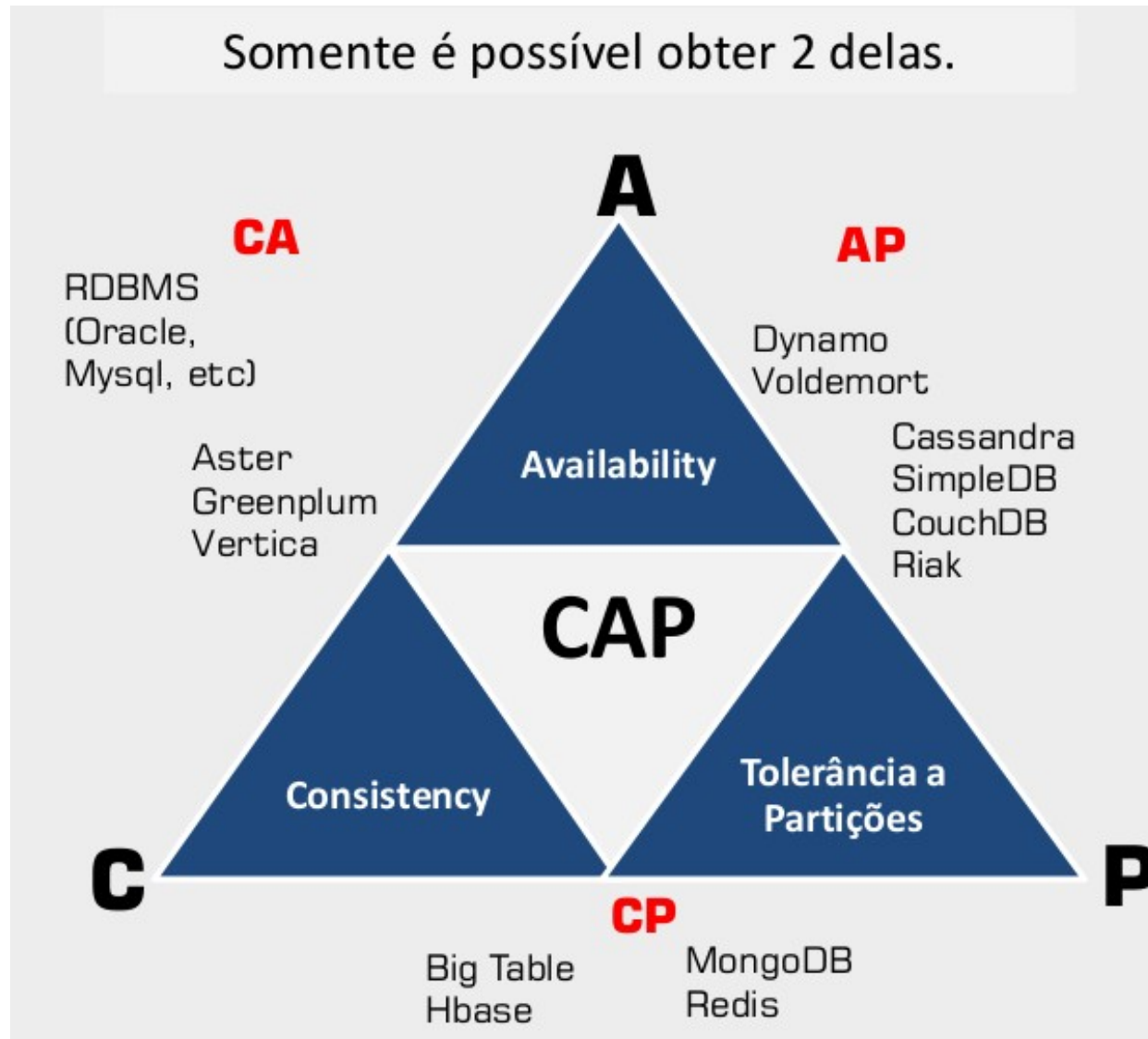
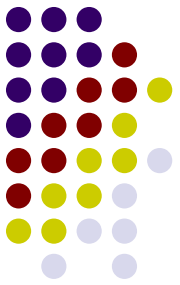
Se em uma escrita a comunicação com algum nó falhar, os dados são replicados apenas aos servidores disponíveis sem que o serviço fique indisponível (escrita ainda é possível). Após o reestabelecimento da comunicação, os dados são sincronizados com os demais servidores.

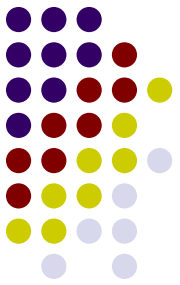


## 6.4 - Sistemas AP



# 7 – Visão Geral CAP





# Referências

- Sadalage, Pramod J.; Fowler, Martin. NOSQL Essencial , Editora Novatec, 2013
- Gilbert, S.; Lynch, N. A. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant web services. SIGACT News, 33(2):51–59, 2002.