

UNIVERSITY OF ILLINOIS  
AT URBANA-CHAMPAIGN

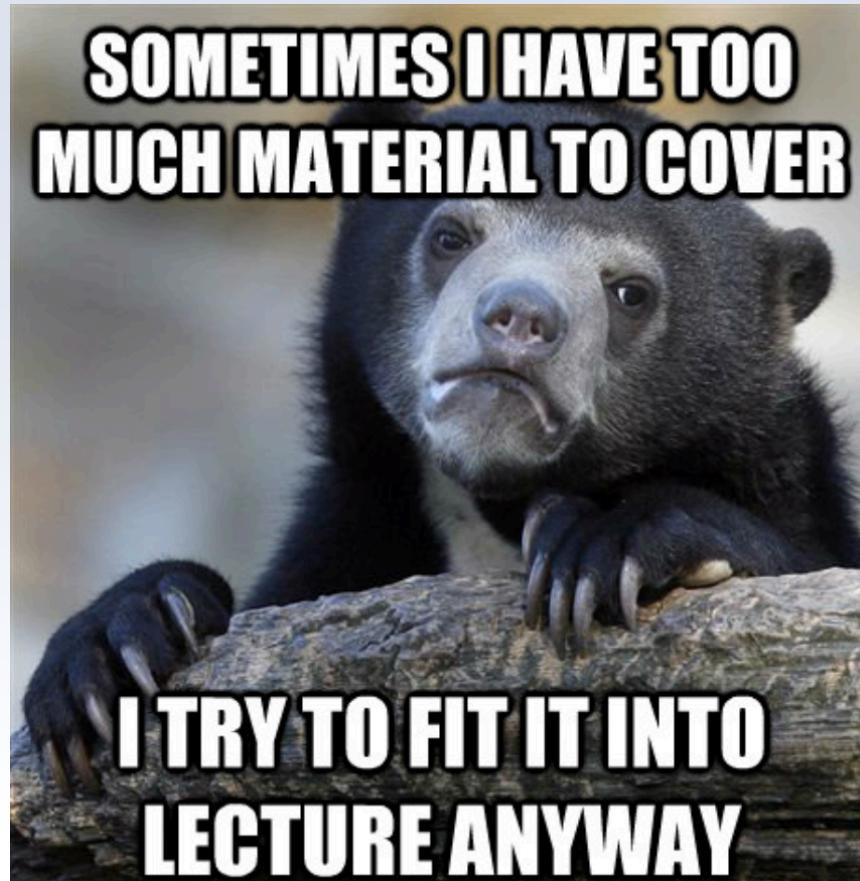
# CS411 - Logging and Recovery



[illinois.edu](http://illinois.edu)



# Mea Culpa



# Cost Estimation

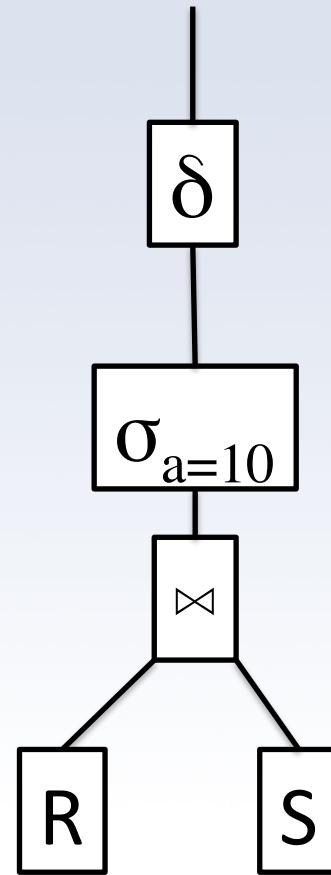
- ***How*** we do it: not that important
  - Can obviously do better
  - Estimation equations provided for exam
- That we ***can*** do it is important



# Example

```
SELECT DISTINCT *
FROM R, S
WHERE R.a=10
```

$\delta(\sigma_{a=10}(R \bowtie S))$



# Example

Given  $R(a,b)$  and  $S(b,c)$

$$T(R)=5000, V(R,a)=50, V(R,b)=100$$

$$T(S)=2000, V(S,b)=200, V(S,c)=100$$

$$U=\delta(\sigma_{a=10}(R \bowtie S))$$

$$T(R \bowtie S)=5000*2000/200=50000$$
$$T(\sigma_{a=10}(R \bowtie S))=50000/50=1000$$

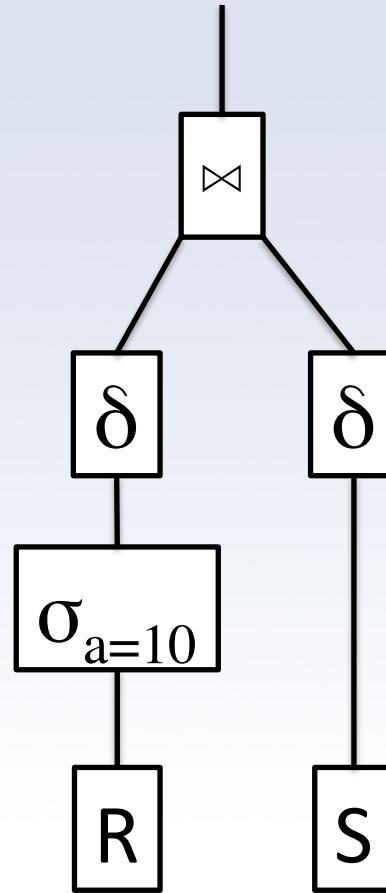
cost=51,000

$$T(U)=\min(100,5000)=1000$$



# Example

```
SELECT DISTINCT *
FROM R, S
WHERE R.a=10
```

$$\delta(\sigma_{a=10}(R)) \bowtie \delta(S)$$


# Example

Given  $R(a,b)$  and  $S(b,c)$

$T(R)=5000, V(R,a)=50, V(R,b)=100$

$T(S)=2000, V(S,b)=200, V(S,c)=100$

$U = \delta(\sigma_{a=10}(R)) \bowtie \delta(S)$

$T(\sigma_{a=10}(R))=100, T(\delta(\sigma_{a=10}(R)))=50$

$T(\delta(S))=1000$

$\text{cost}=1,150$

$T(U)=1000*50/200=250$



# Example

```
SELECT a,b,AVG(x)
FROM R,S,T,U,V,W
WHERE C
GROUP BY a,b
```



# Naïve Approach

$$\gamma_{a,b,\text{avg}(x)}(\sigma_C(R \bowtie S \bowtie T \bowtie U \bowtie V \bowtie W))$$

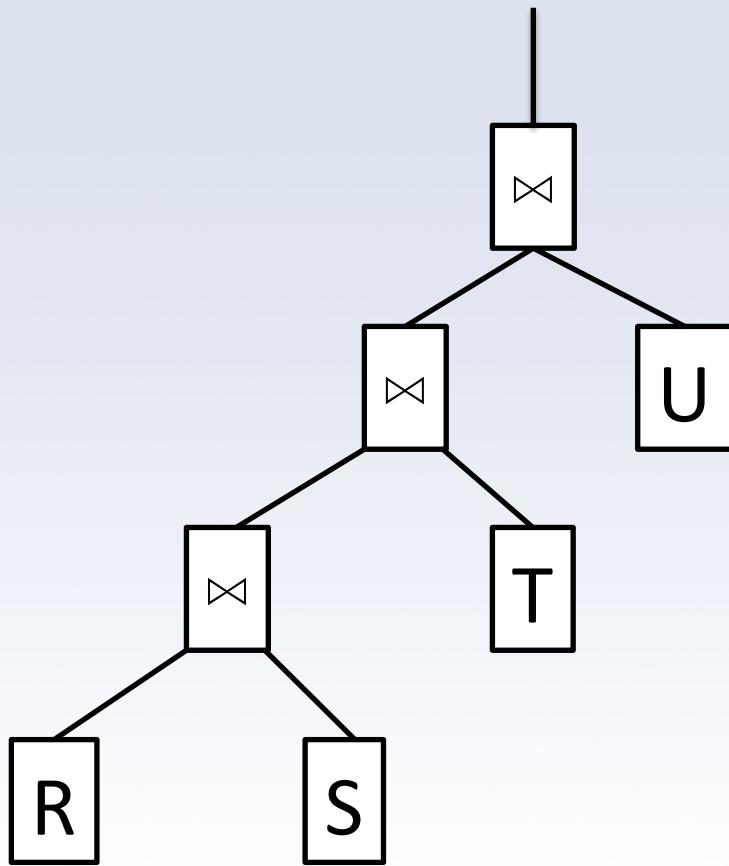
- Number of binary trees with  $n$  leaves =  $(n-1)$  Catalan number =  $[2(n-1)]! / [(n!)(n-1)!]$
- Number of arrangements of  $n$  leaves =  $n!$
- Number of ways to reorder  $n$  joins =  $[2(n-1)]! / (n-1)!$

Reordering JUST the joins in this problem: 30,240

For 10 joins: 17,643,225,600



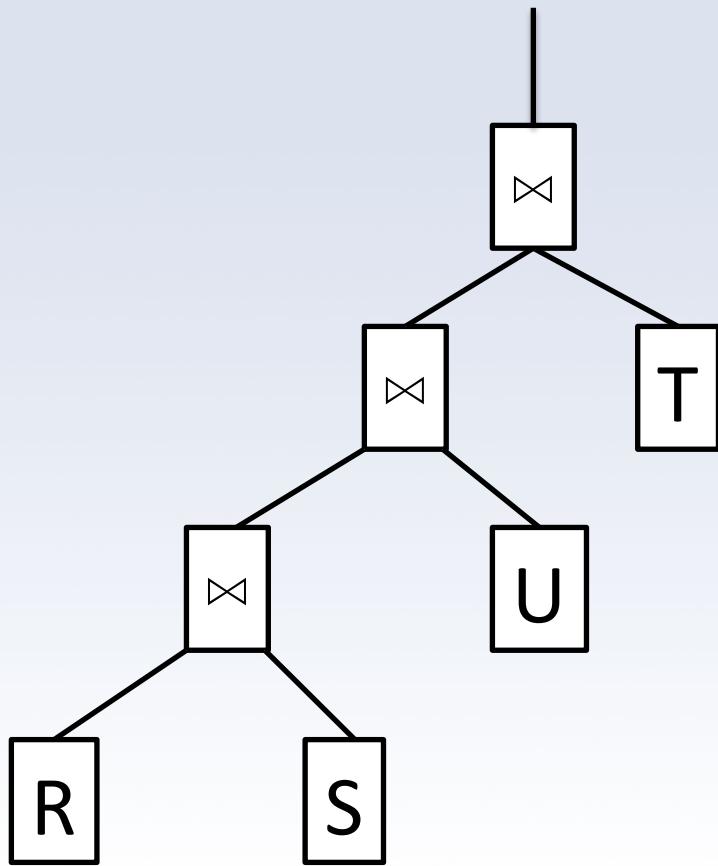
# Left Join Tree



- A tree topology with two advantages:
  1. Only  $n!$  different orderings possible
  2. Can be easily pipelined with our algorithms
- On examples/exams, only use left join trees



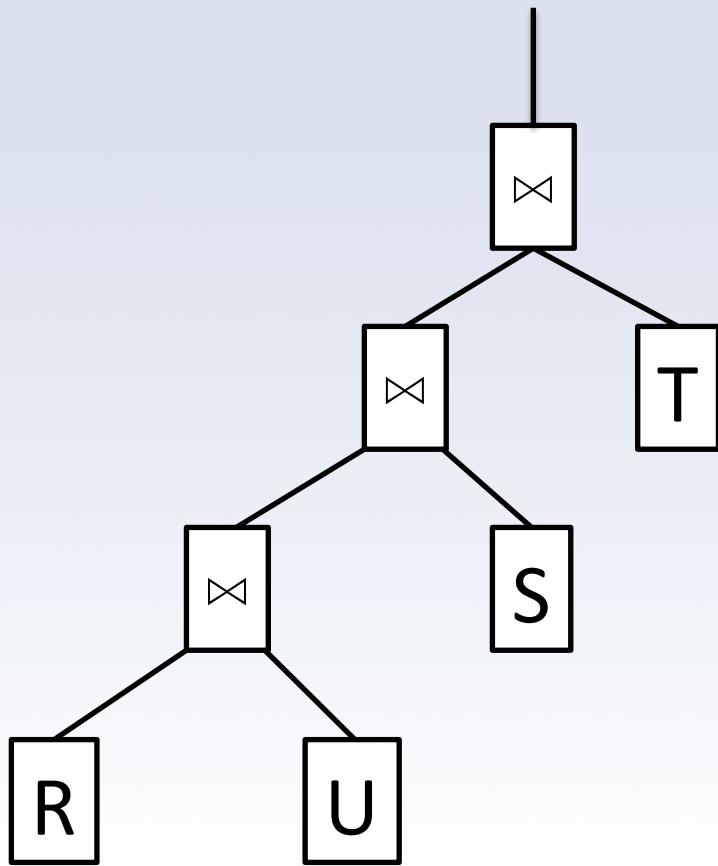
# Left Join Tree



- A tree topology with two advantages:
  1. Only  $n!$  different orderings possible
  2. Can be easily pipelined with our algorithms
- On examples/exams, only use left join trees



# Left Join Tree



- A tree topology with two advantages:
  1. Only  $n!$  different orderings possible
  2. Can be easily pipelined with our algorithms
- On examples/exams, only use left join trees



# Dynamic Programming

- Main idea: in large combinatorial problems, many of the partial results are helpful in later steps
- Example:  
 $\text{Cost}((R \bowtie S) \bowtie T)$  result is useful for computing  
 $\text{Cost}((R \bowtie S) \bowtie T \bowtie U)$   
 $\text{Cost}(((R \bowtie S) \bowtie T) \bowtie V)$   
 $\text{Cost}(((R \bowtie S) \bowtie T) \bowtie W)$

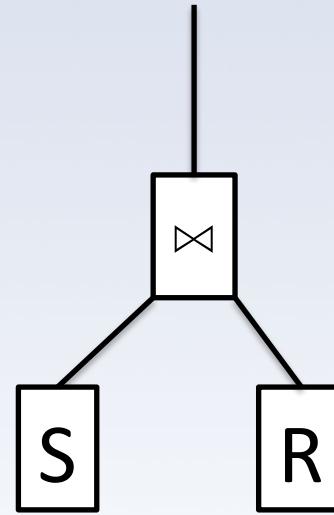
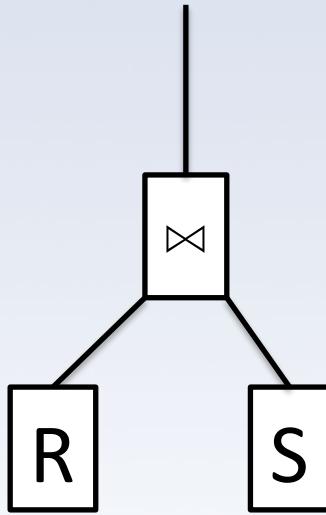


# Example

Subset	Size(S)	Cost	Plan
{R,S}			
{R,T}			
{R,U}			
{S,T}			
{S,U}			
{T,U}			
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Subsets of size 2



# Example

Subset	Size(S)	Cost	Plan
{R,S}			$R \bowtie S$
{R,T}			
{R,U}			
{S,T}			
{S,U}			
{T,U}			
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	$T(R)T(S)/\max(V(R,b),V(S,b))$		$R \bowtie S$
{R,T}			
{R,U}			
{S,T}			
{S,U}			
{T,U}			
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	1k*1k/200		$R \bowtie S$
{R,T}			
{R,U}			
{S,T}			
{S,U}			
{T,U}			
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}			
{R,U}			
{S,T}			
{S,U}			
{T,U}			
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			

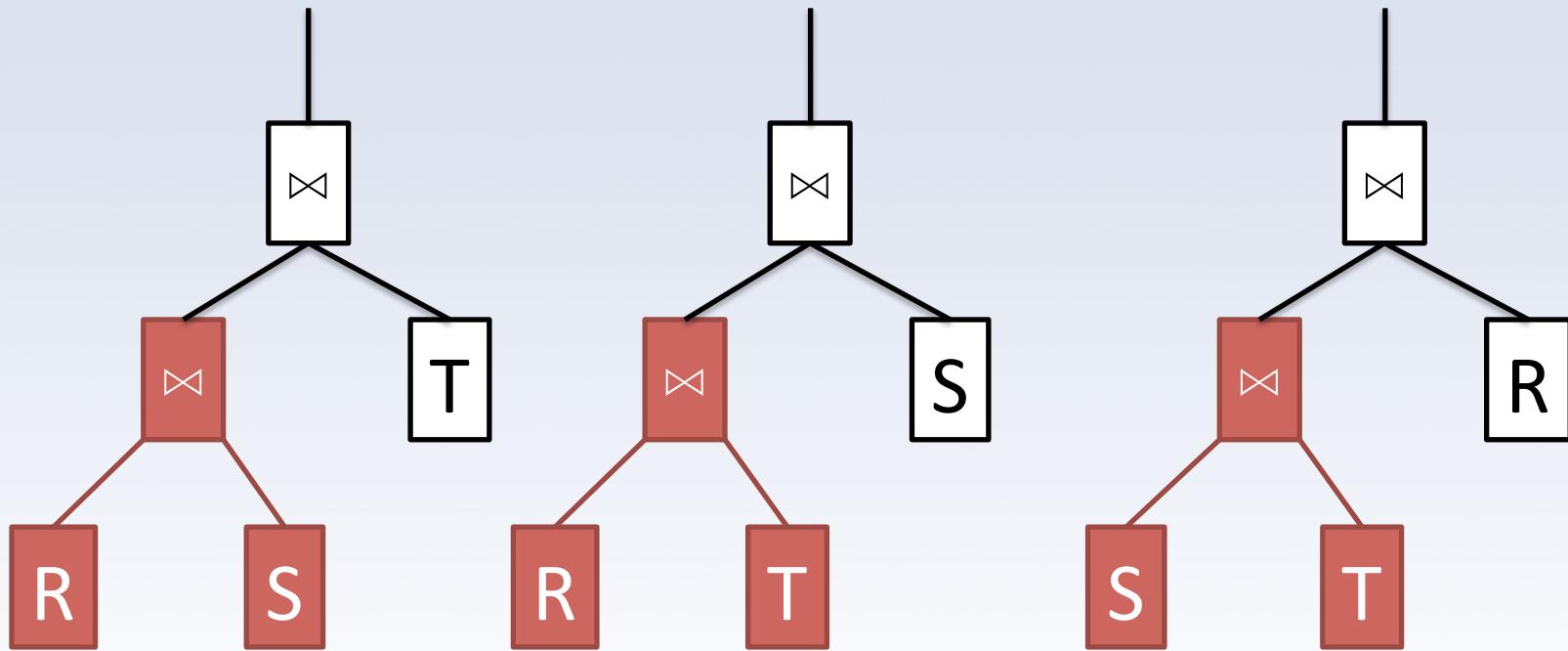


# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}			
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Subsets of size 3



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	$(1k * 1k * 1k) / (500 * 200) = 10k$	5k+0	$(R \bowtie S) \bowtie T$
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	$(1k * 1k * 1k) / (500 * 200) = 10k$	1,000k+0	$(R \bowtie T) \bowtie S$
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	$(1k * 1k * 1k) / (500 * 200) = 10k$	2k+0	$(S \bowtie T) \bowtie R$
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}			
{R,T,U}			
{S,T,U}			
{R,S,T,U}			

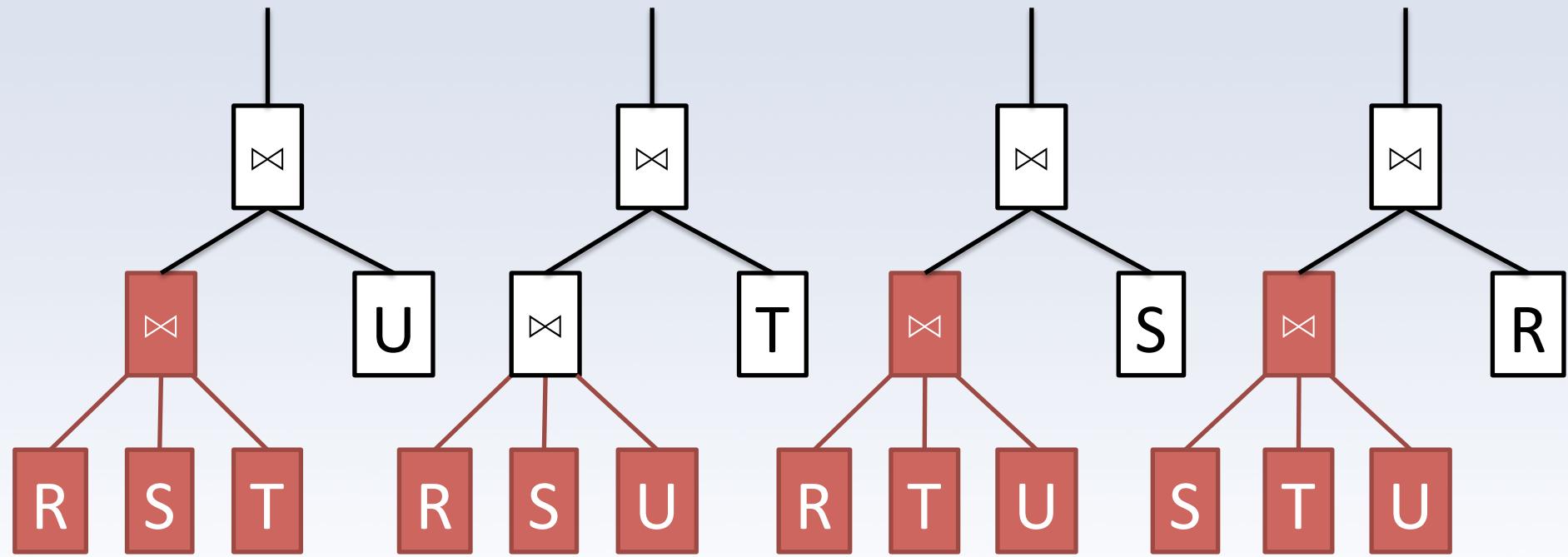


# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}			



# Subsets of size 4



# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}	$(1k * 1k * 1k * 1k) / (1k * 500 * 200 * 100) = 10k$	$10k + 2k + 0 + 0 = 12k$	$((S \bowtie T) \bowtie R) \bowtie U$

# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}	$(1k * 1k * 1k * 1k) / (1k * 500 * 200 * 100) = 10k$	$50k + 5k + 0 + 0 = 55k$	$((R \bowtie S) \bowtie U) \bowtie T$

# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}	$(1k * 1k * 1k * 1k) / (1k * 500 * 200 * 100) = 10k$	$10k + 1k + 0 + 0 = 11k$	$((T \bowtie U) \bowtie R) \bowtie S$

# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}	$(1k * 1k * 1k * 1k) / (1k * 500 * 200 * 100) = 10k$	$2k + 1k + 0 + 0 = 3k$	$((T \bowtie U) \bowtie S) \bowtie R$

# Example

Subset	Size(S)	Cost	Plan
{R,S}	5k	0	$R \bowtie S$
{R,T}	1,000k	0	$R \bowtie T$
{R,U}	10k	0	$R \bowtie U$
{S,T}	2k	0	$S \bowtie T$
{S,U}	1,000k	0	$S \bowtie U$
{T,U}	1k	0	$T \bowtie U$
{R,S,T}	10k	2k	$(S \bowtie T) \bowtie R$
{R,S,U}	50k	5k	$(R \bowtie S) \bowtie U$
{R,T,U}	10k	1k	$(T \bowtie U) \bowtie R$
{S,T,U}	2k	1k	$(T \bowtie U) \bowtie S$
{R,S,T,U}	10k	3k	$((T \bowtie U) \bowtie S) \bowtie R$



# Physical Query Plan

- At each node of logical query plan:
  - Replace relational algebra operator
  - Substitute execution algorithms
    - How much memory is available?
    - Are relations indexed?
    - Are results sorted?
    - Should we pipeline or materialize intermediate results?

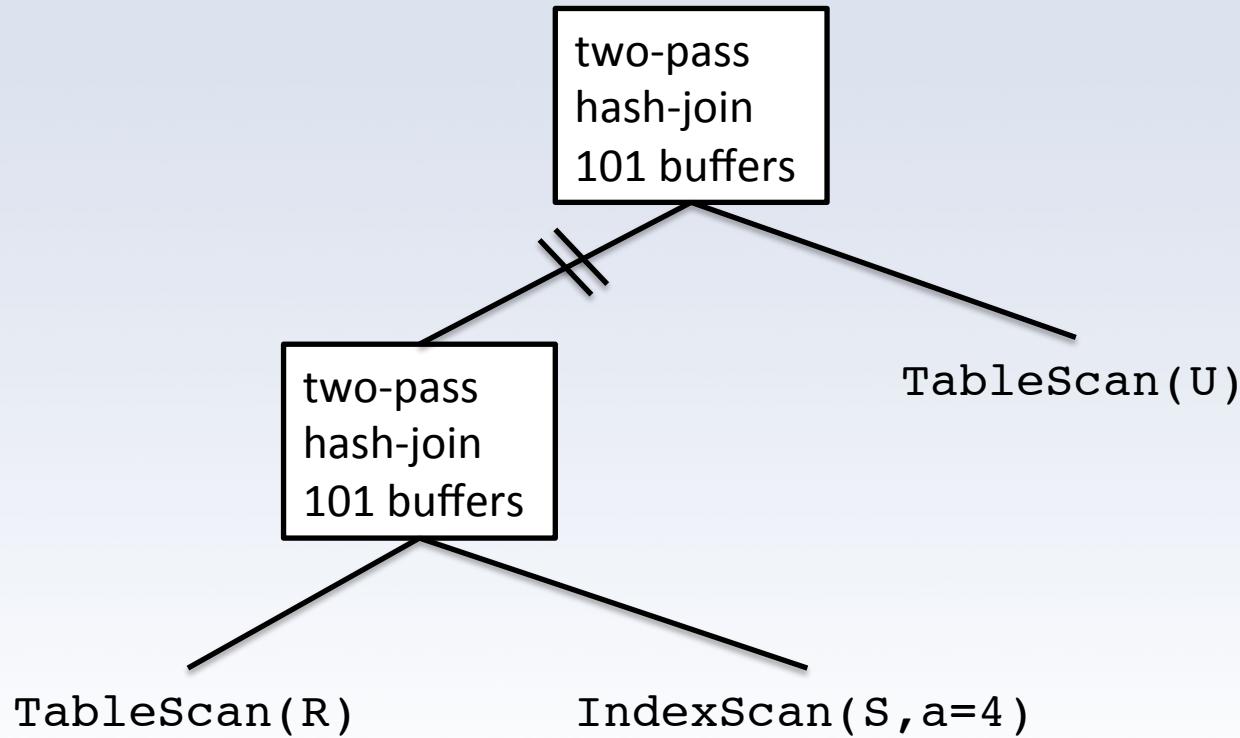


# Pipelining vs. Materialization

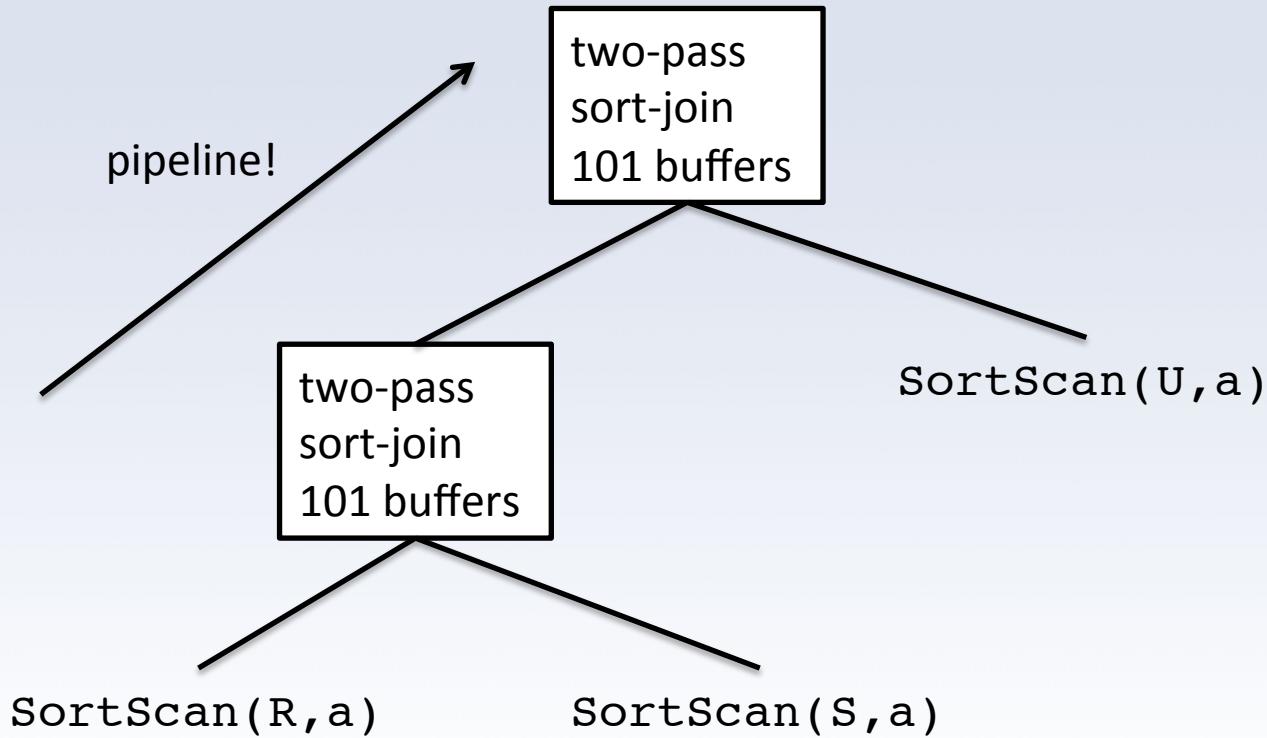
- Materialization - write the results of the operation back to the disk
  - results in 2 intermediate disk I/Os per block
- Pipelining - keep results in memory buffers
  - Results of one iterator passed to the next



# Physical Query Plan



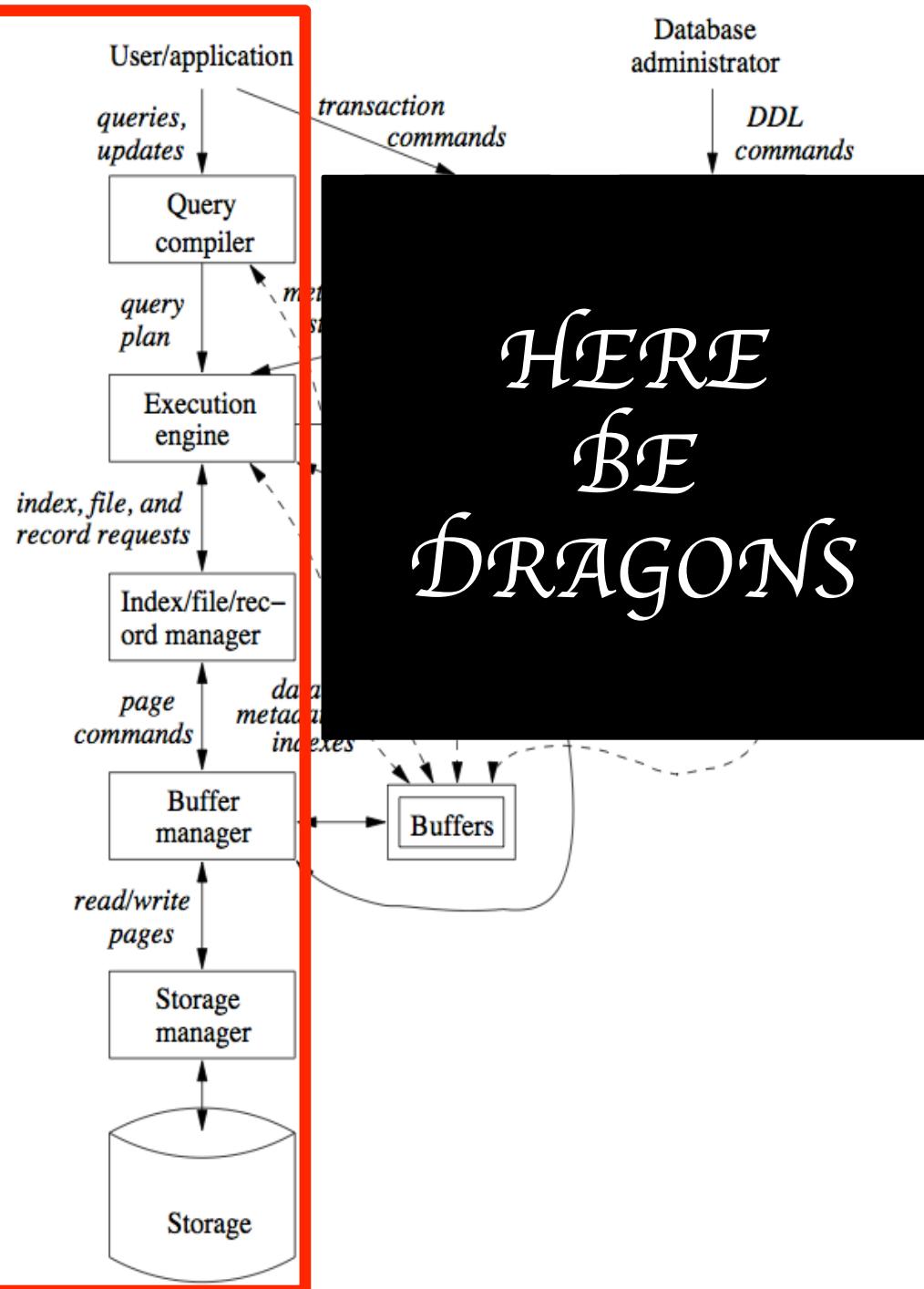
# Physical Query Plan



# MP3

- We implemented a few algorithms for operators (e.g. select, join, project)
- We implemented cost estimates
- We have five queries
  - Improve the query plan for all five
  - Keeping within memory constraint





# User Submits Query

```
SELECT DISTINCT *
FROM R, S
WHERE R.a>3
```



User/application

*queries,  
updates*

Query  
compiler

*query  
plan*

*trans*

*meta  
statis*

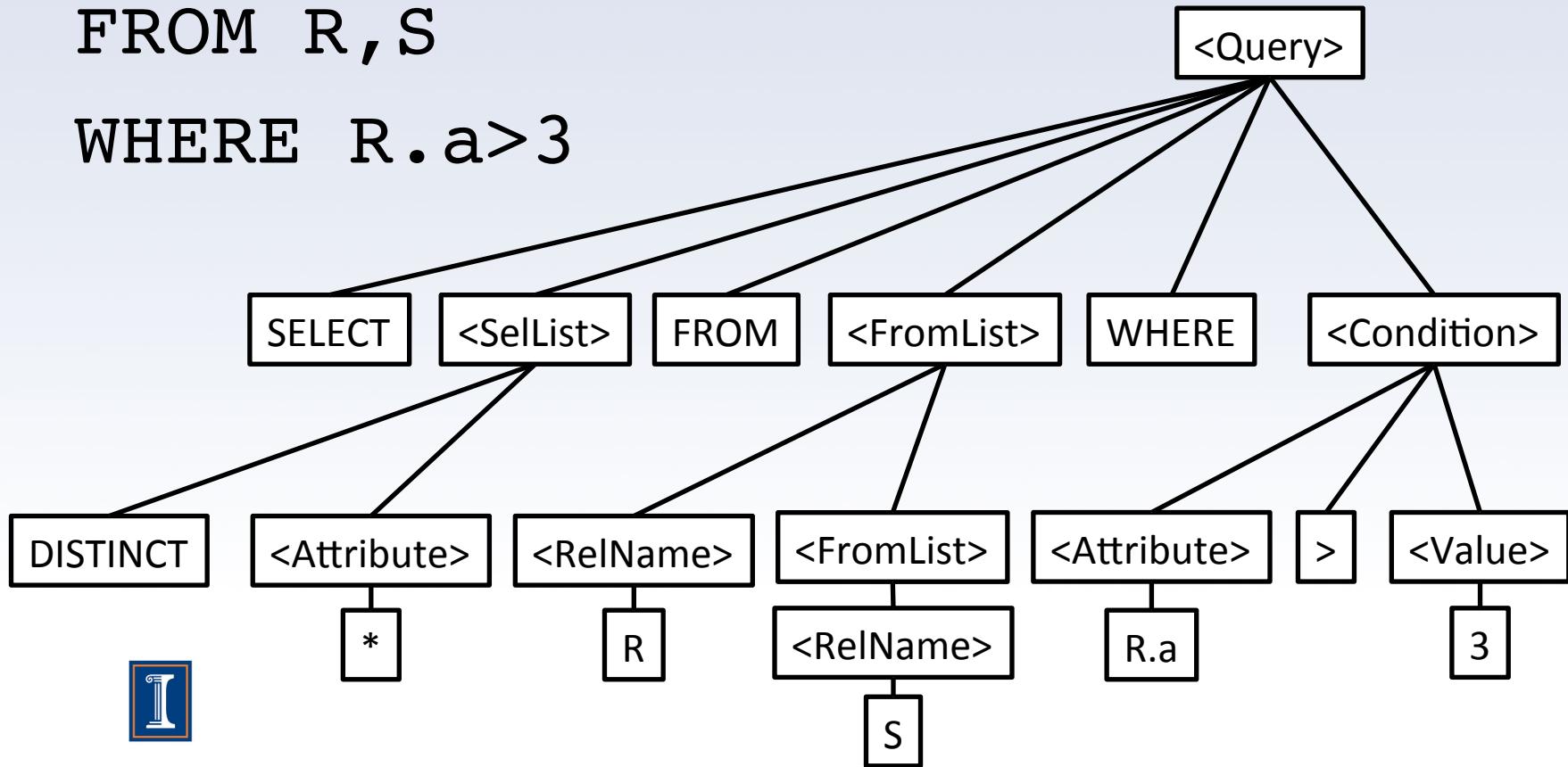


# Parse Tree

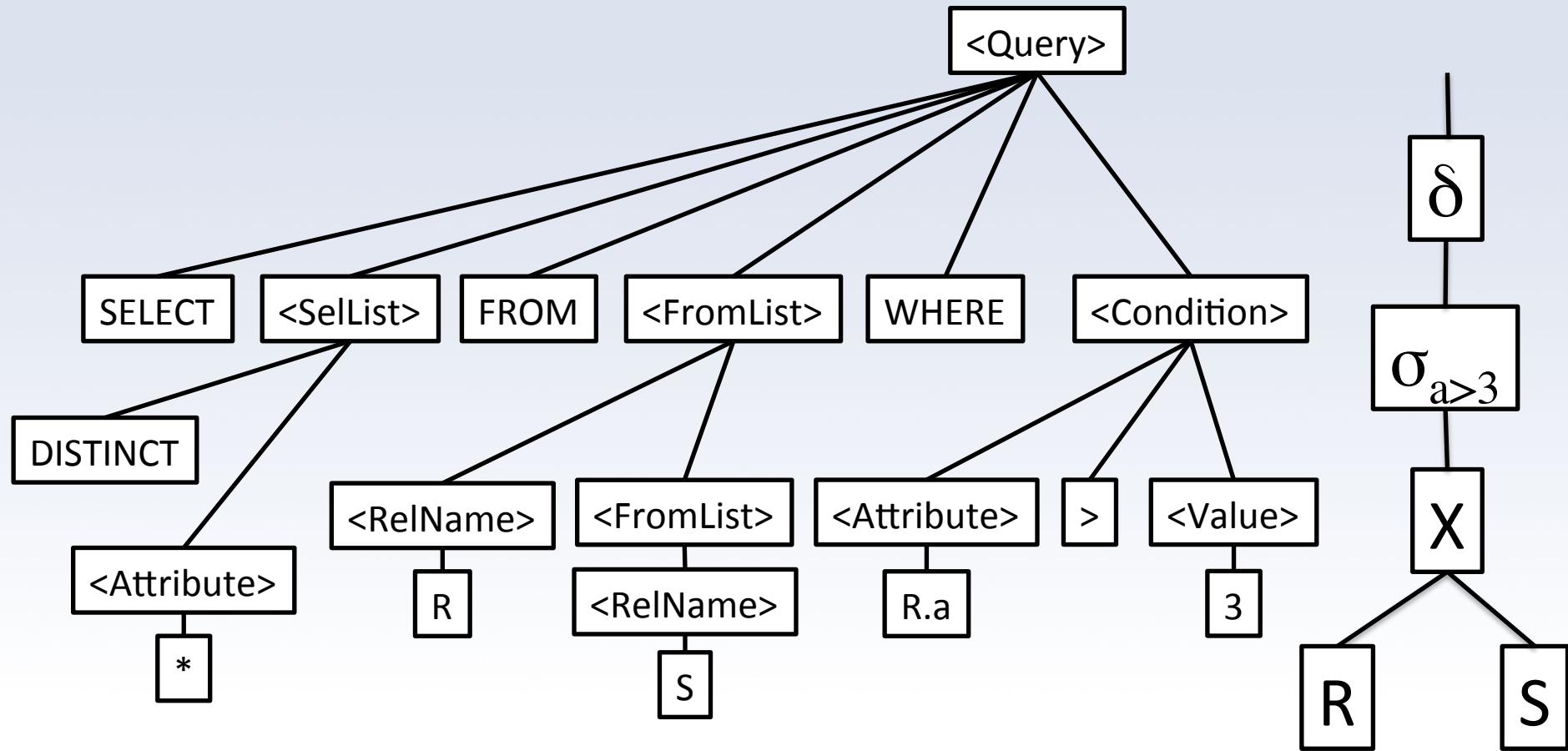
SELECT DISTINCT \*

FROM R, S

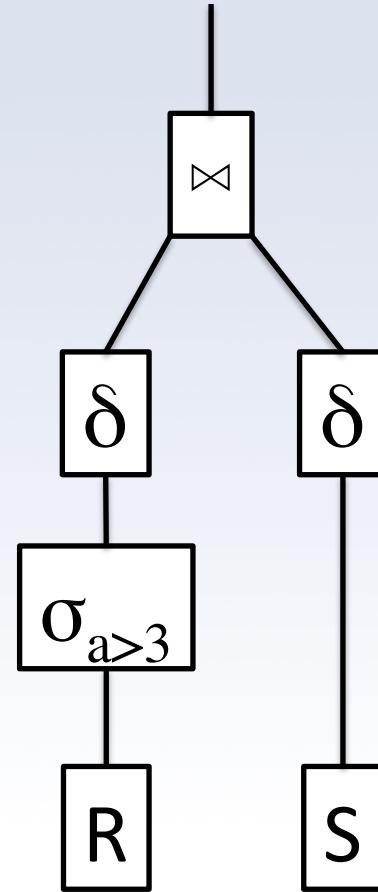
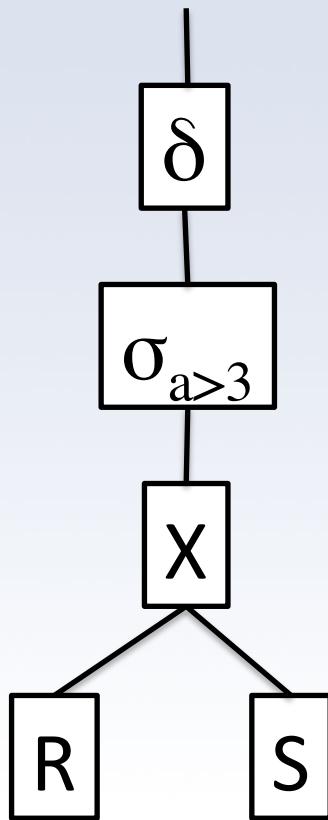
WHERE R.a > 3



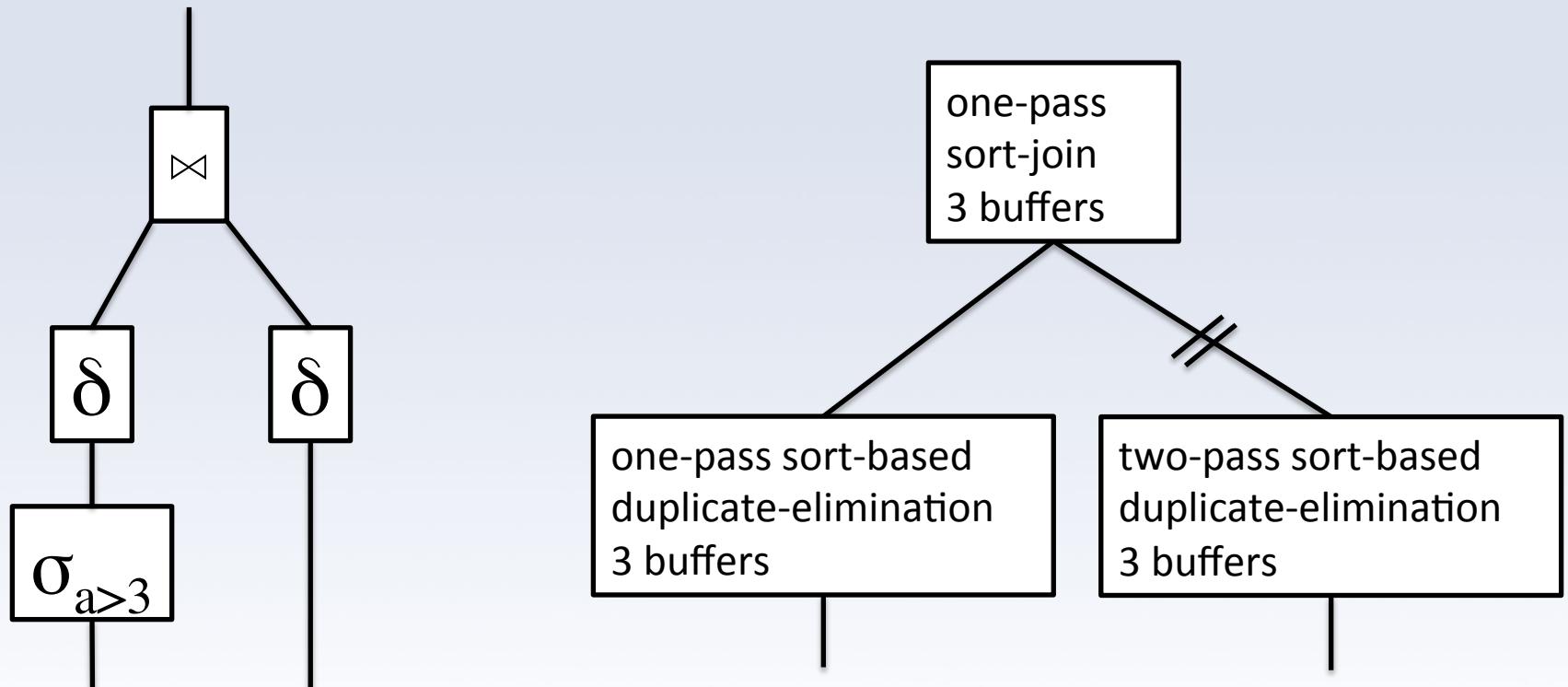
# Logical Query Plan

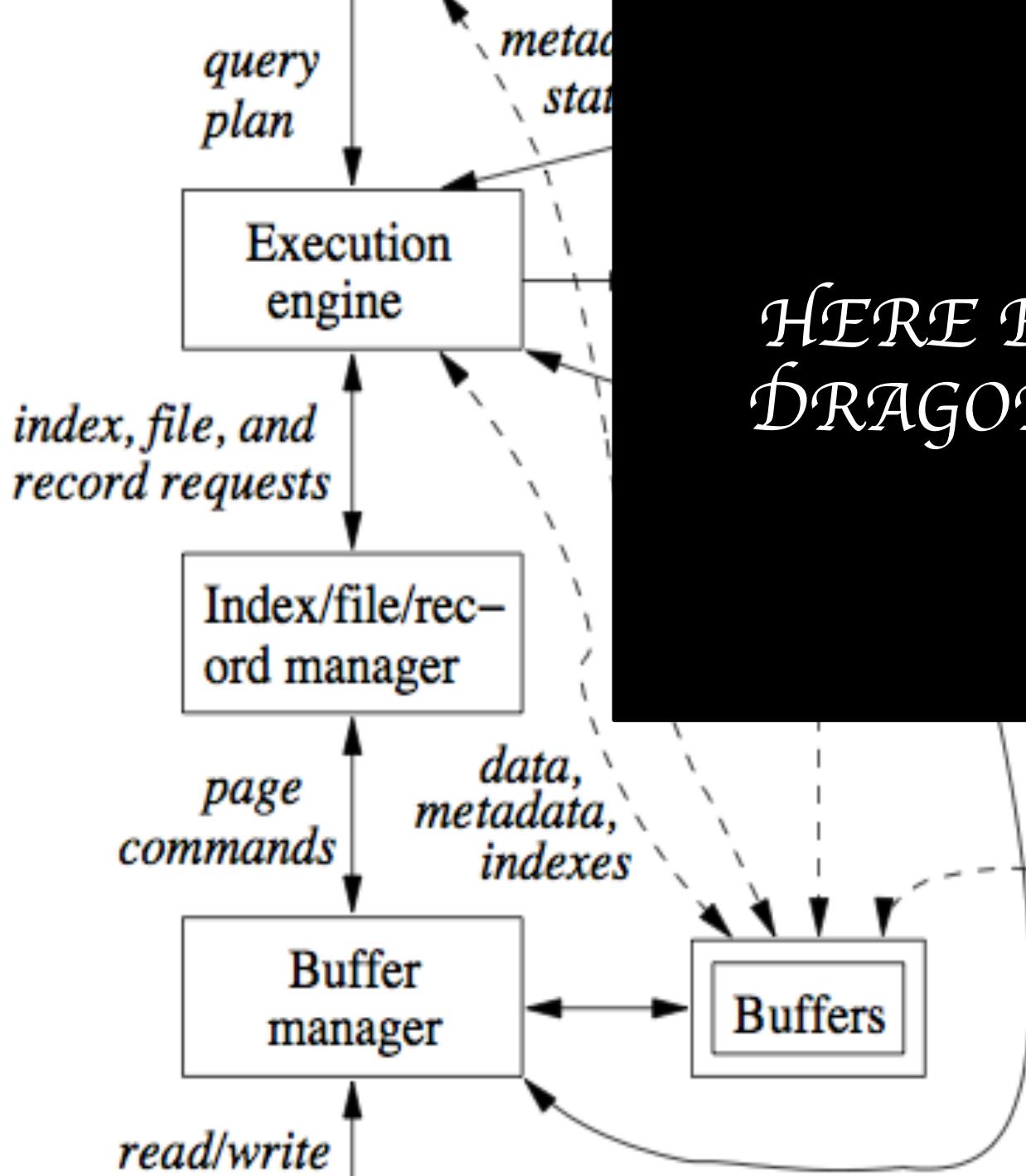


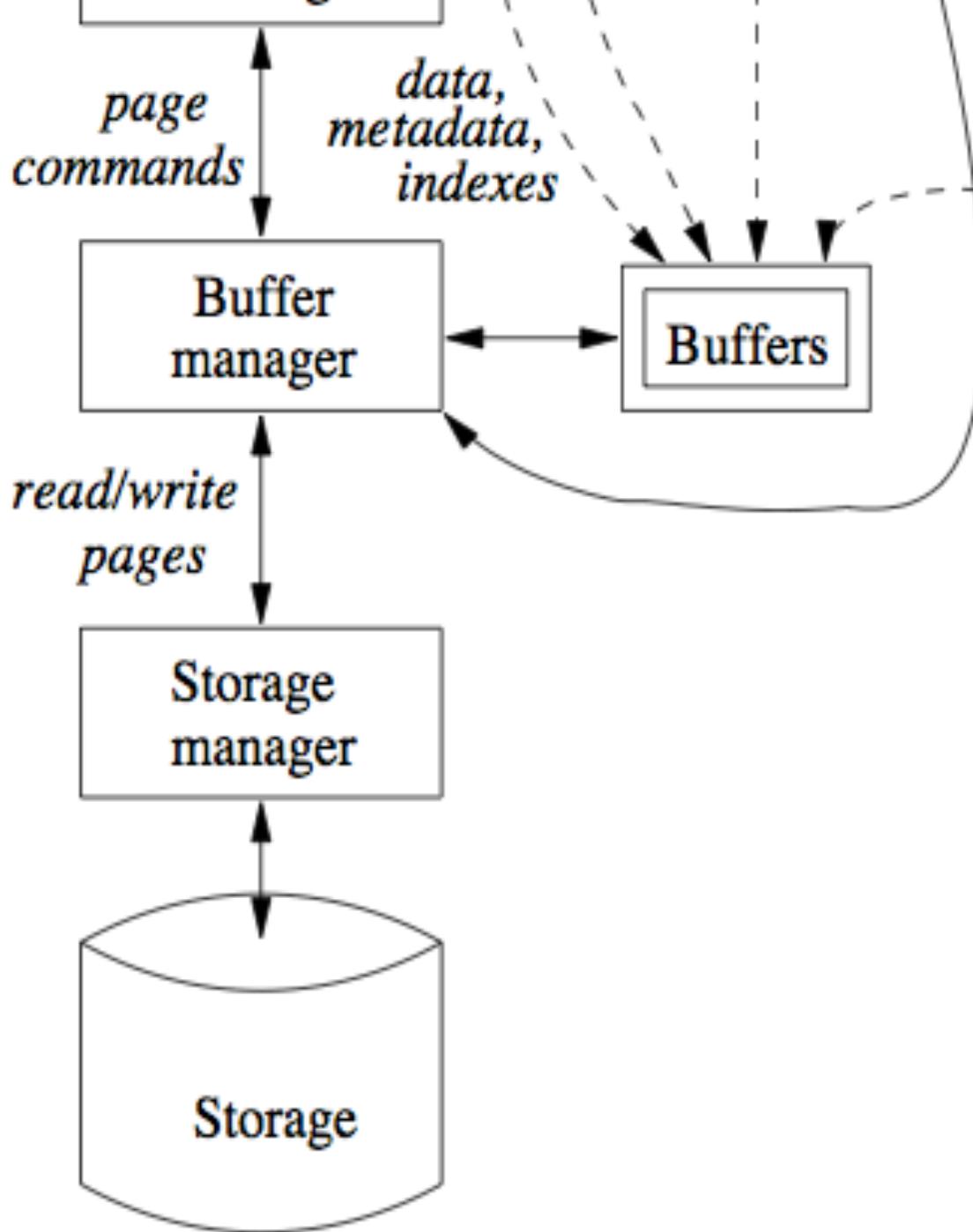
# Logical Query Plan (Optimized)

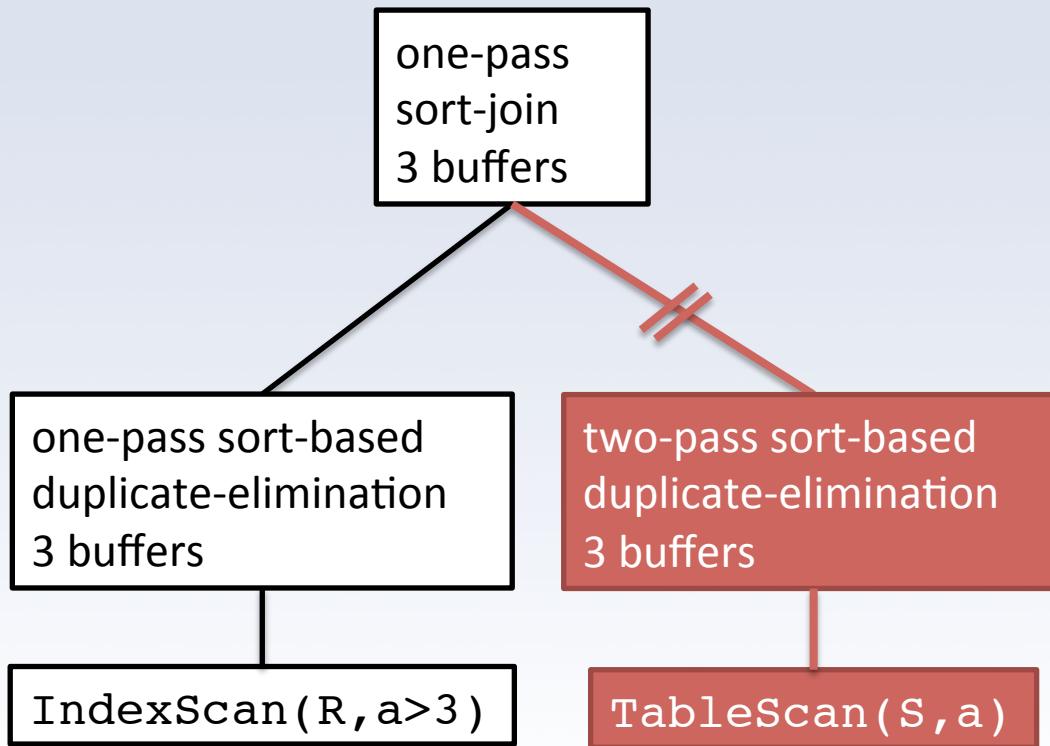


# Physical Query Plan

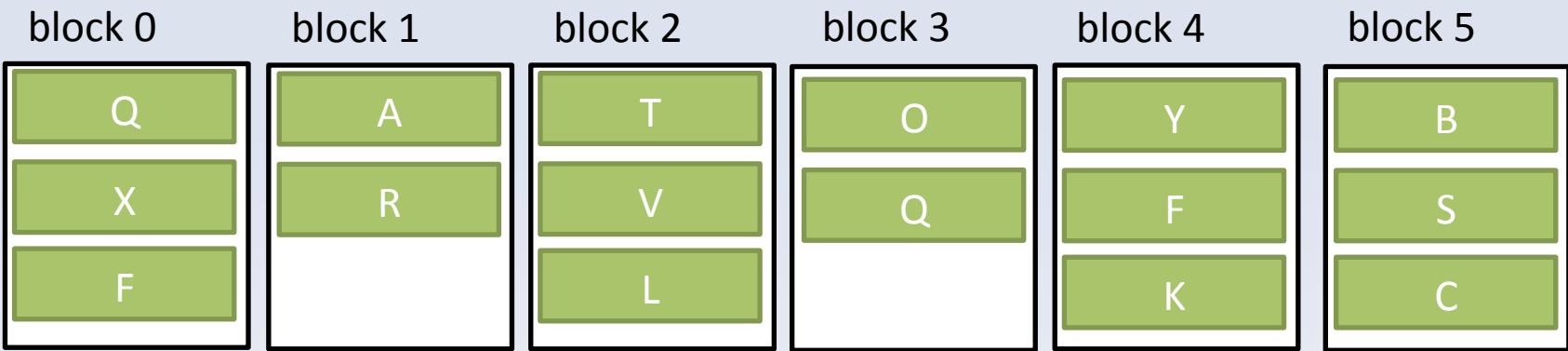






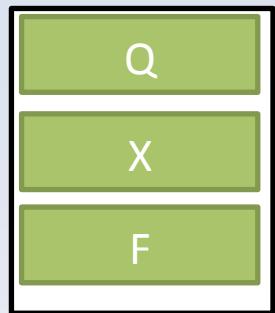


# 2 Pass Duplicate-Elimination

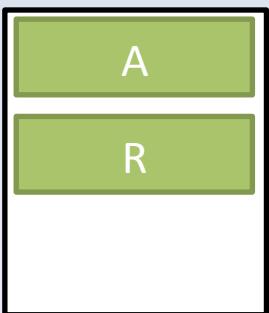


# Phase 1

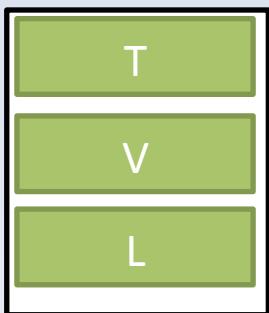
block 0



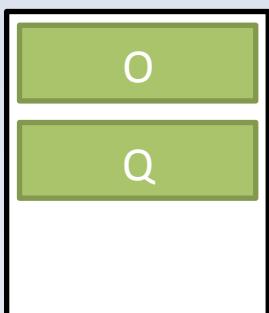
block 1



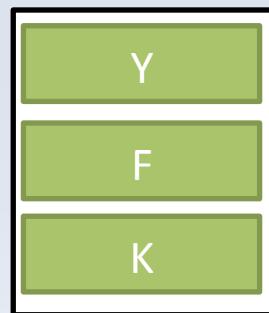
block 2



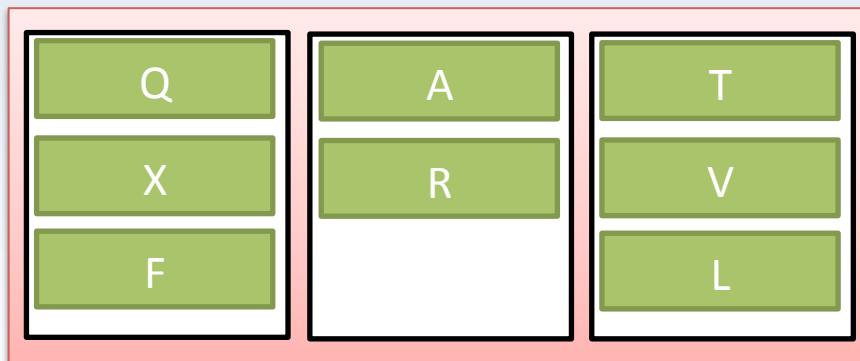
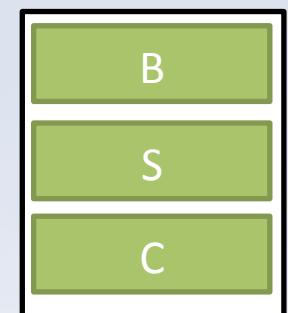
block 3



block 4

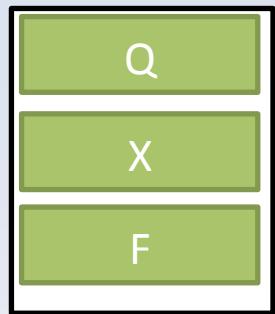


block 5

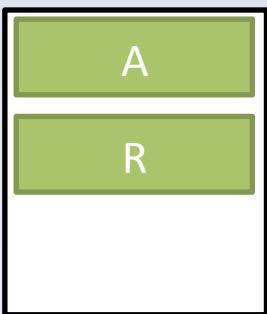


# Phase 1

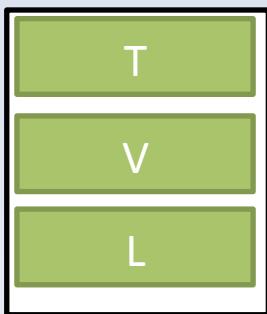
block 0



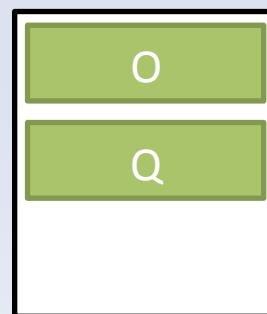
block 1



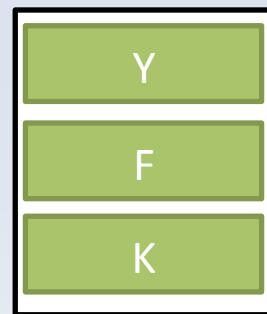
block 2



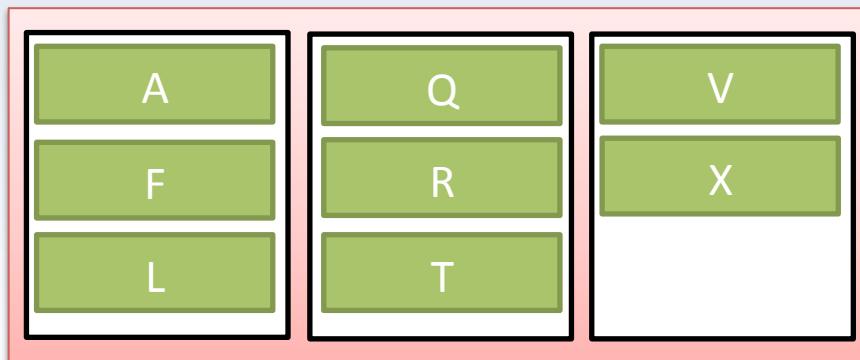
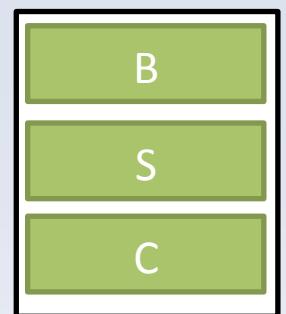
block 3



block 4

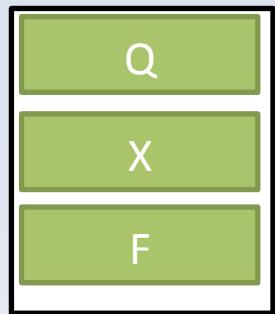


block 5

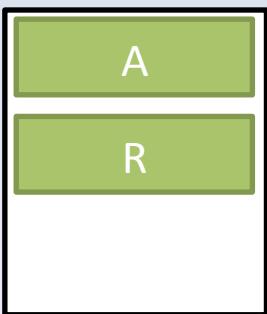


# Phase 1

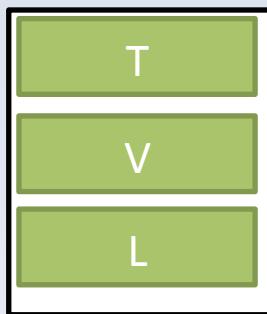
block 0



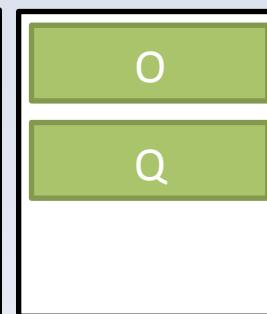
block 1



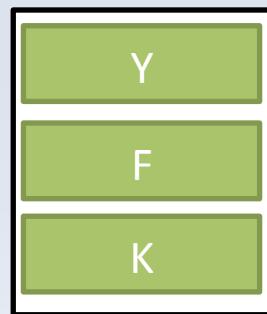
block 2



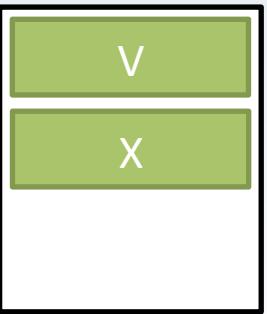
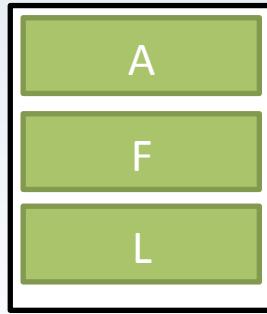
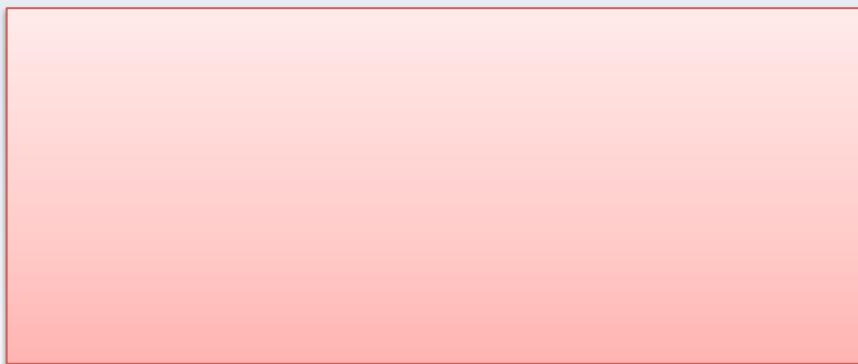
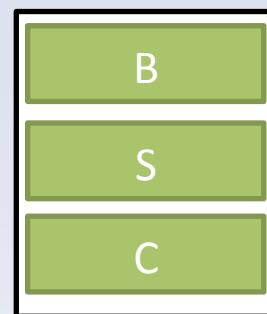
block 3



block 4

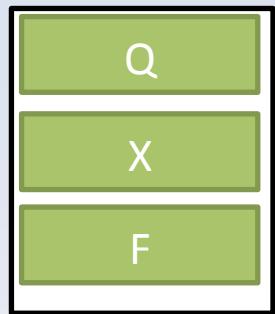


block 5

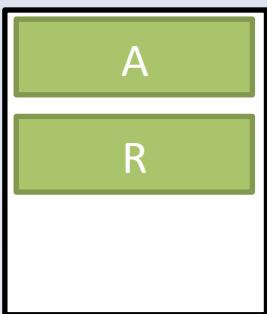


# Phase 1

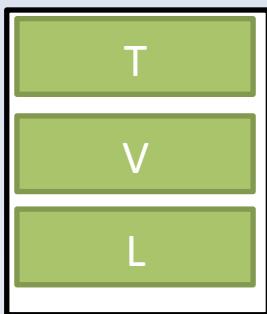
block 0



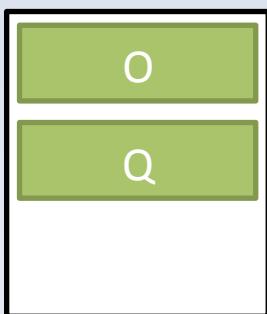
block 1



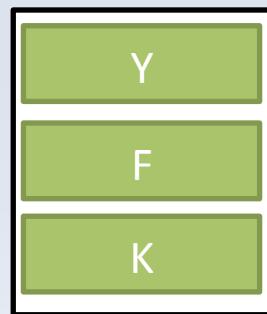
block 2



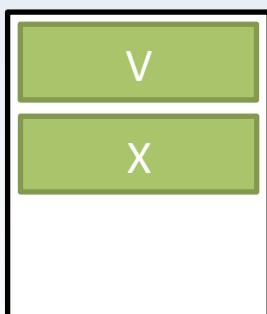
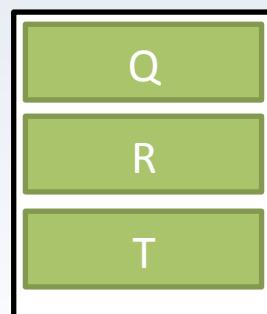
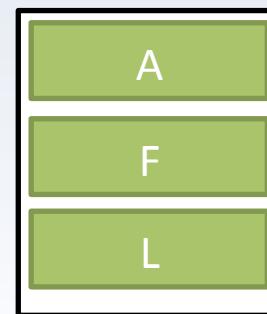
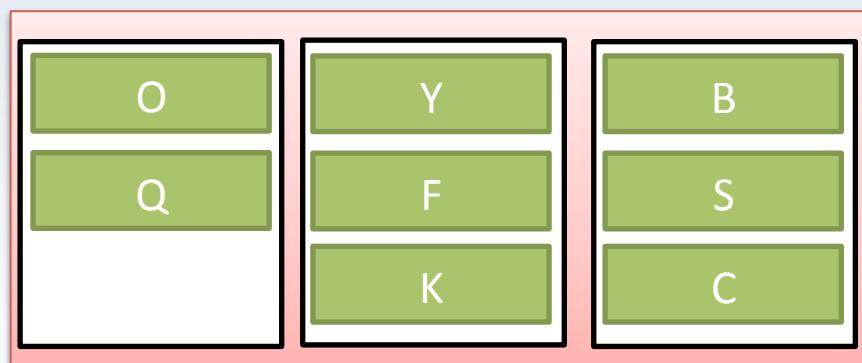
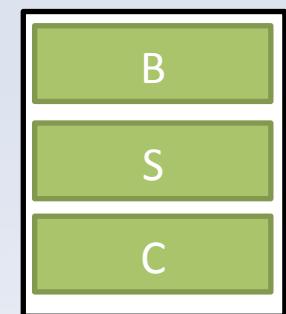
block 3



block 4

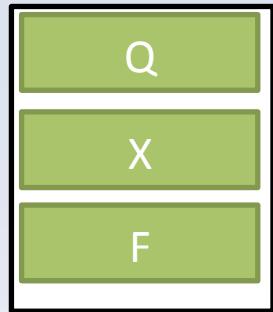


block 5

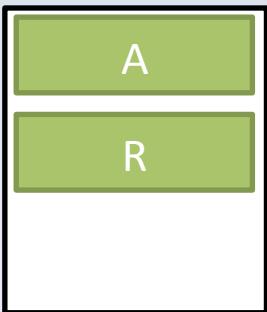


# Phase 1

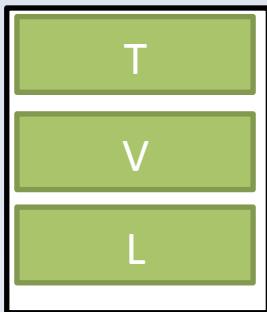
block 0



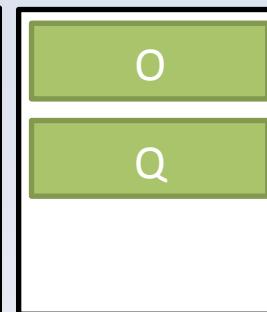
block 1



block 2



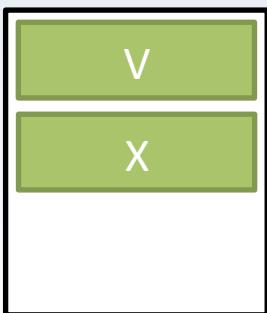
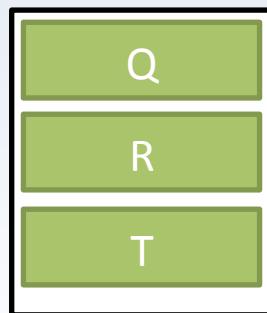
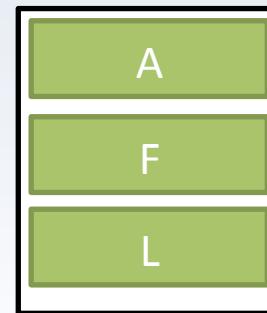
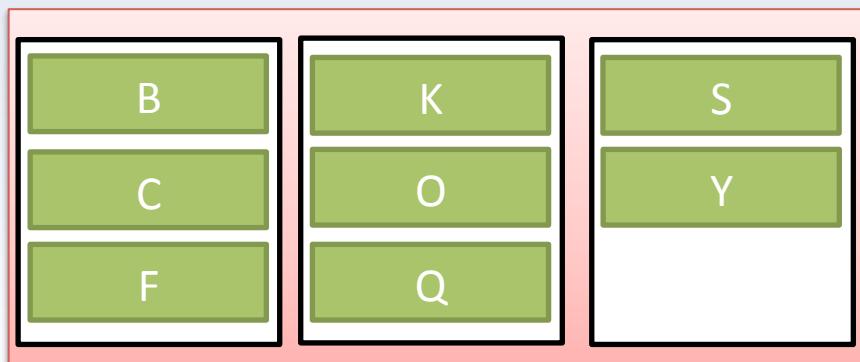
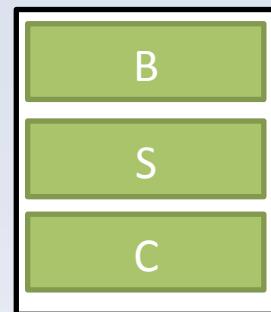
block 3



block 4

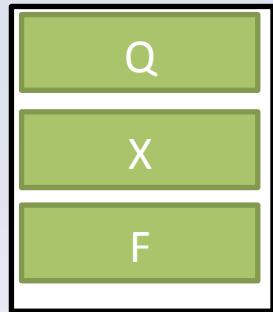


block 5

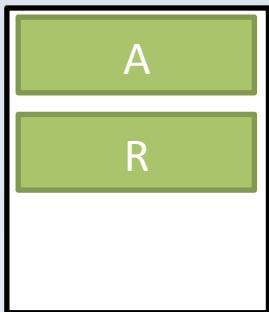


# Phase 1

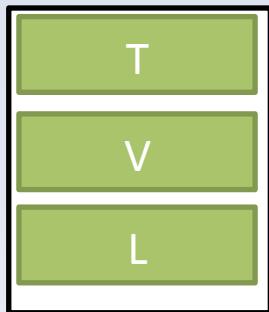
block 0



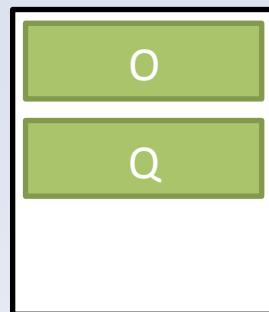
block 1



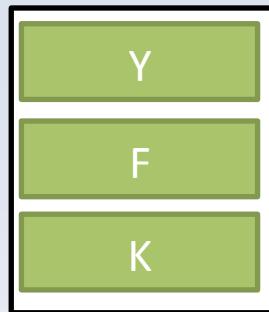
block 2



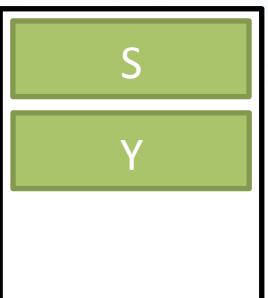
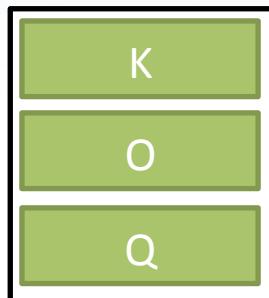
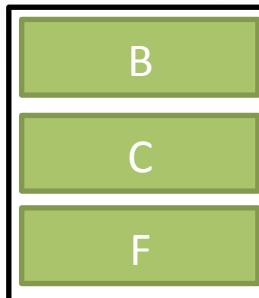
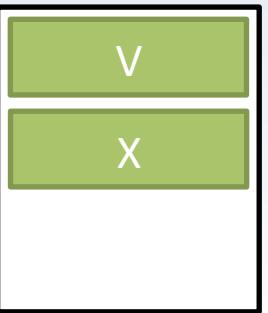
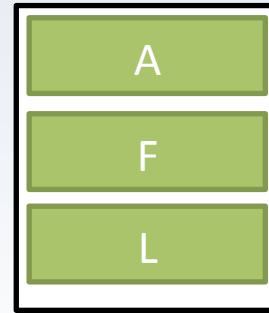
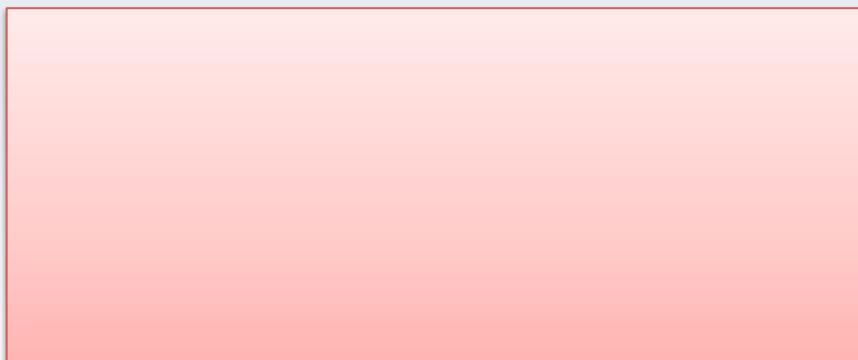
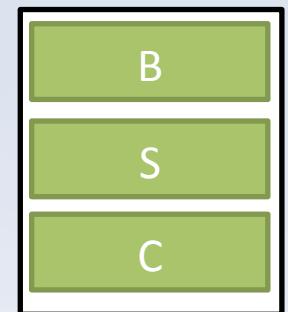
block 3



block 4

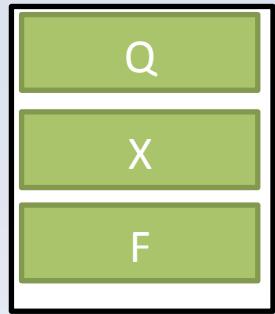


block 5

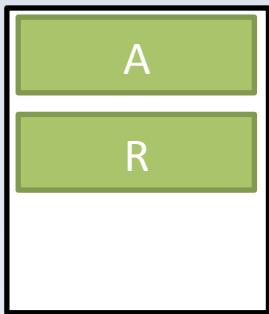


# Phase 2

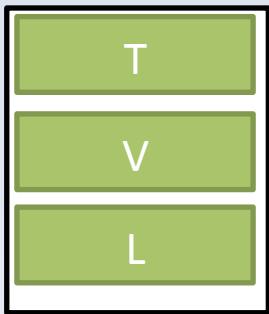
block 0



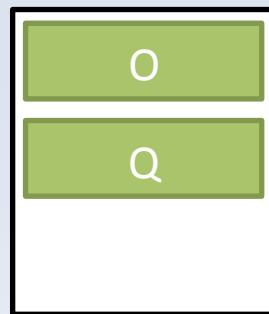
block 1



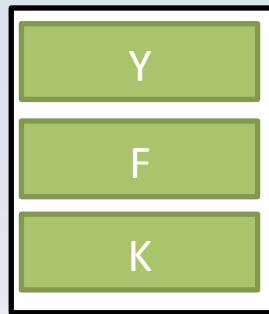
block 2



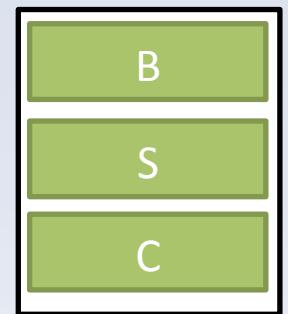
block 3



block 4



block 5



sublist 1

sublist 2

output

A

F

L

Q

R

T

V

X

B

C

F

K

O

Q

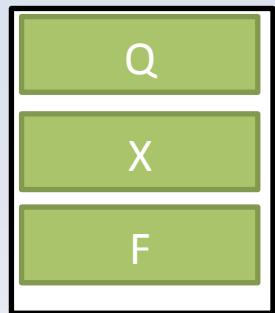
S

Y

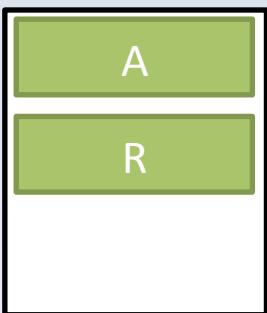


# Phase 2

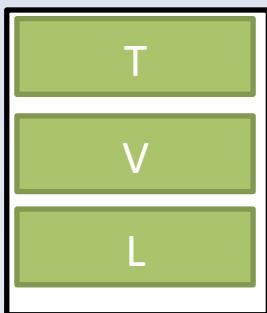
block 0



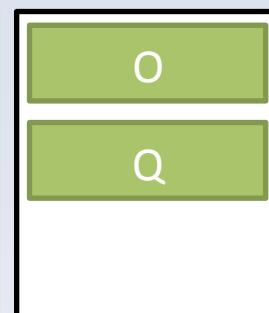
block 1



block 2



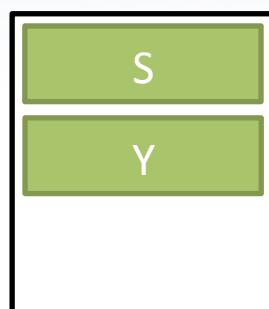
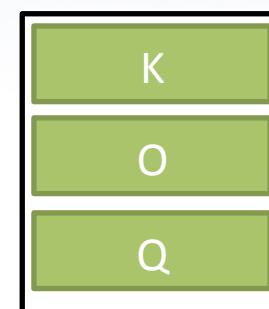
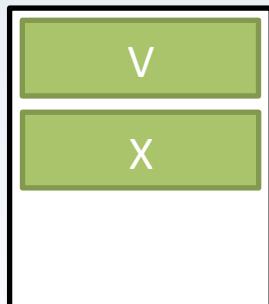
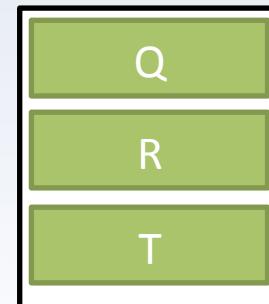
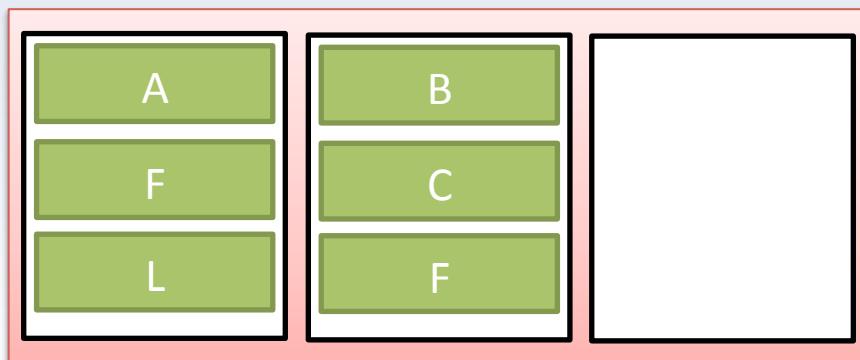
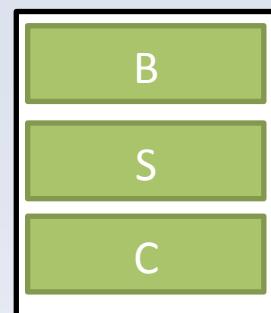
block 3



block 4

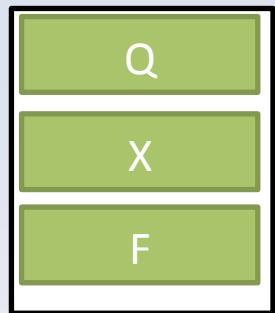


block 5

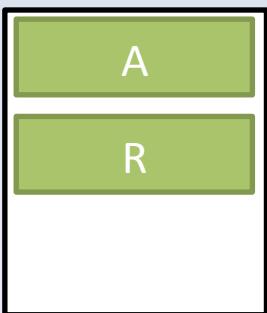


# Phase 2

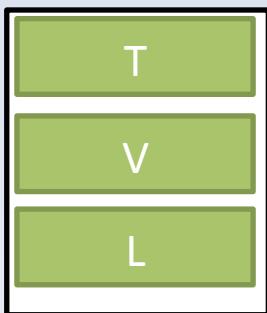
block 0



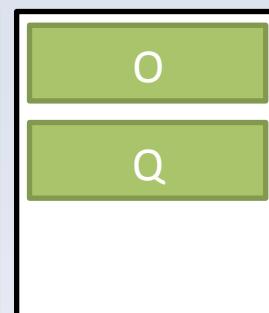
block 1



block 2



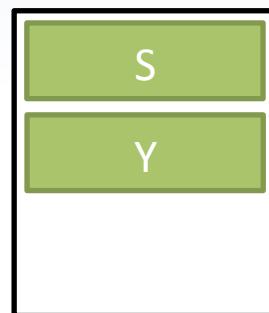
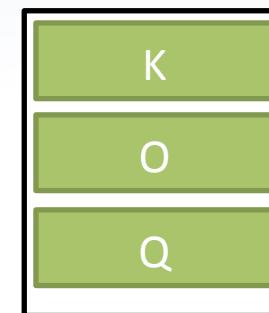
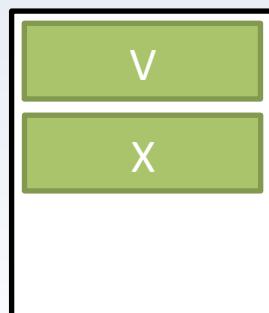
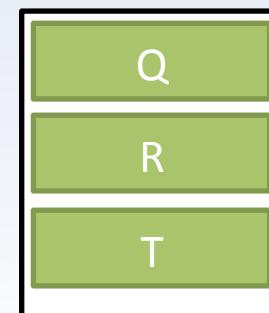
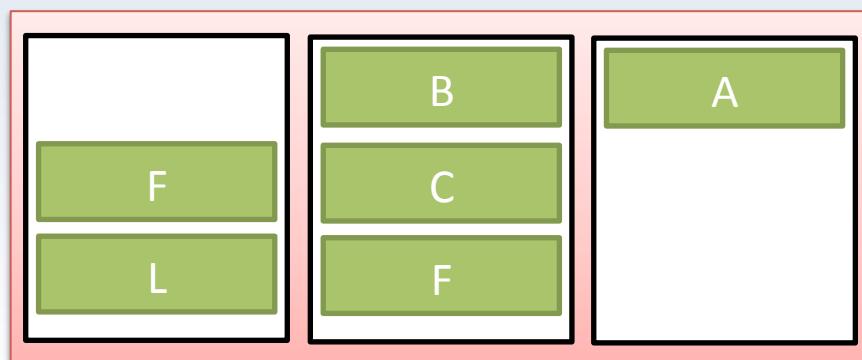
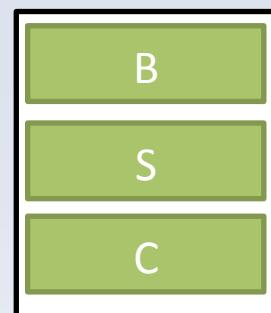
block 3



block 4

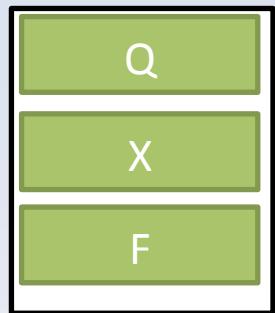


block 5

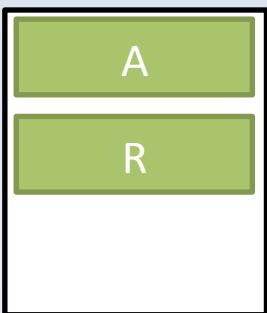


# Phase 2

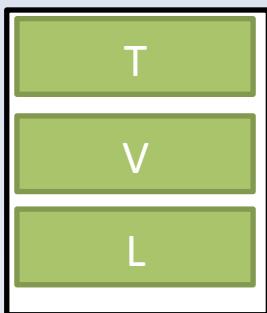
block 0



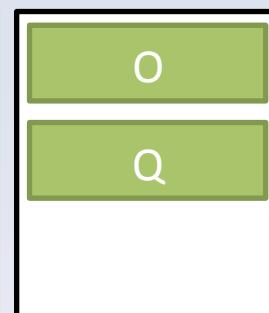
block 1



block 2



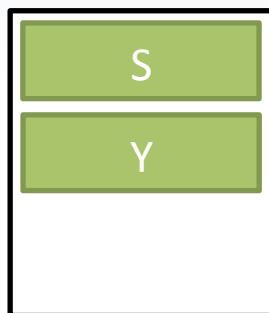
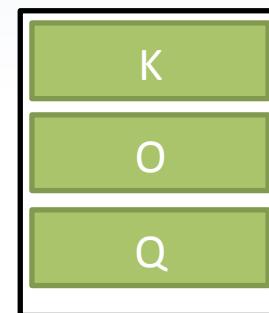
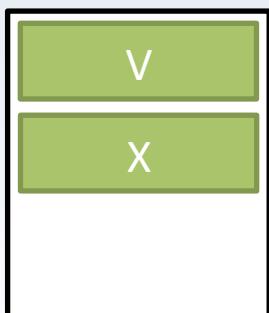
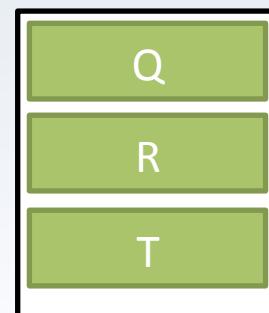
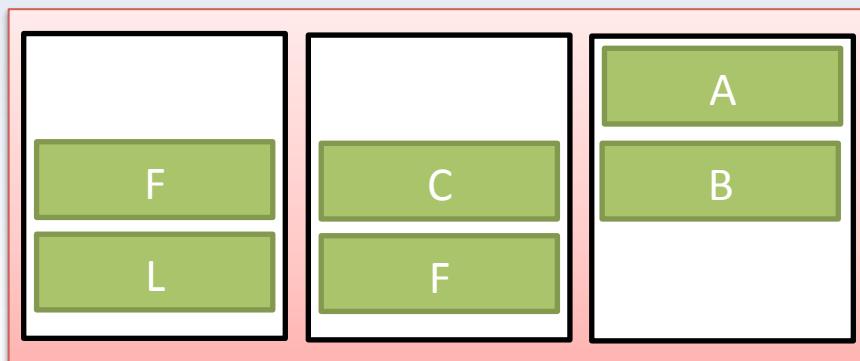
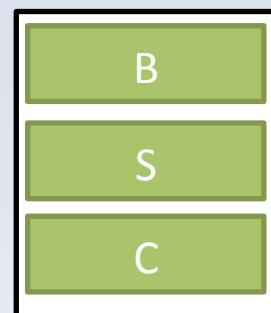
block 3



block 4

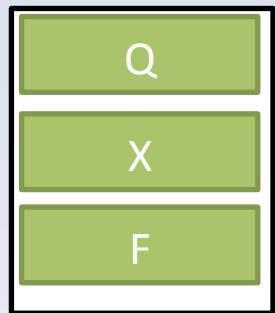


block 5

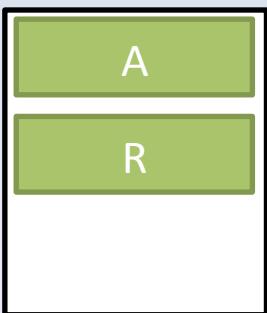


# Phase 2

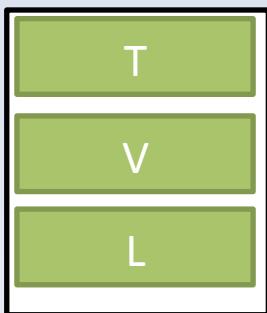
block 0



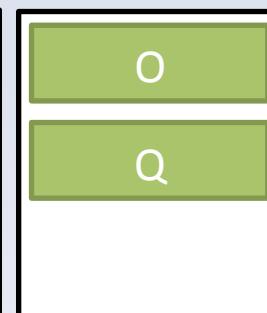
block 1



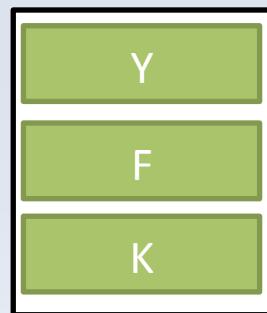
block 2



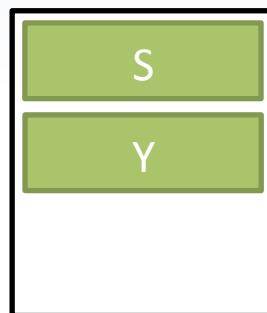
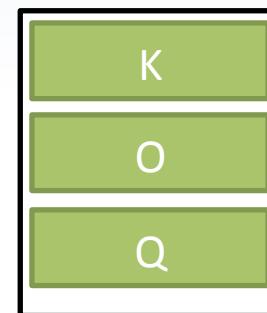
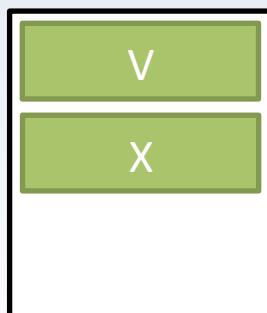
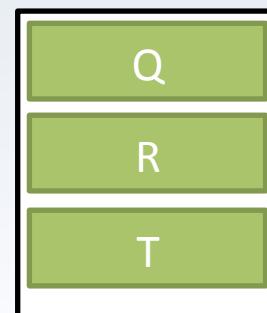
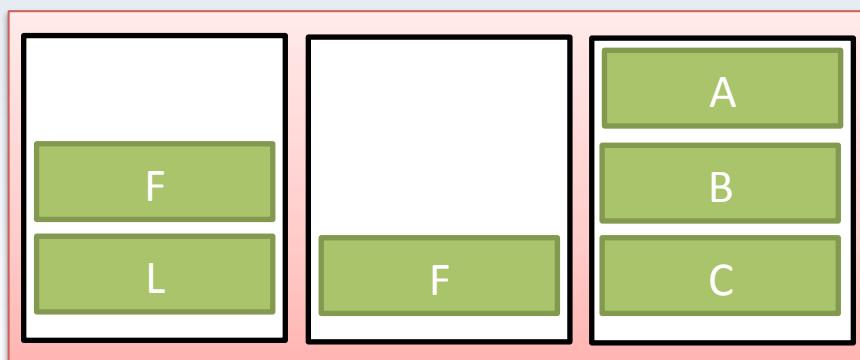
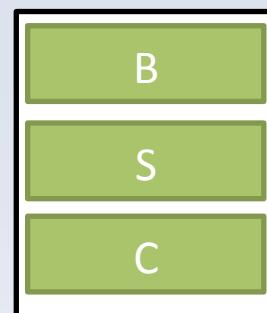
block 3



block 4

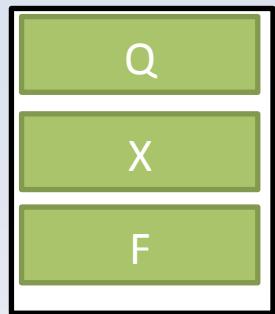


block 5

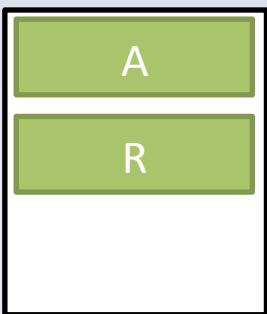


# Phase 2

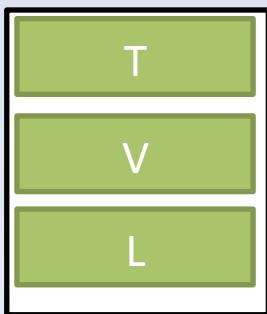
block 0



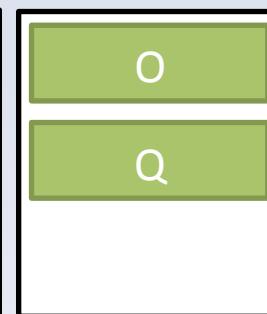
block 1



block 2



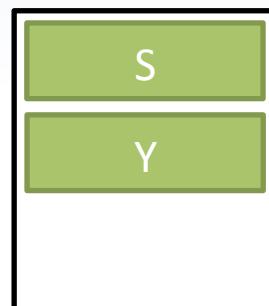
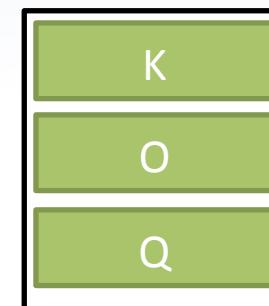
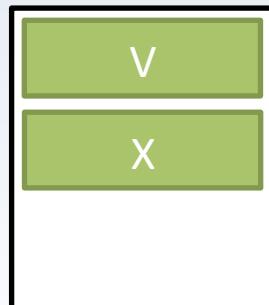
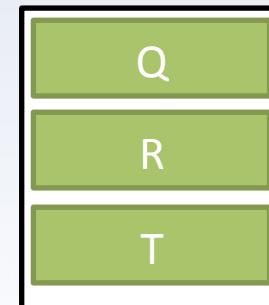
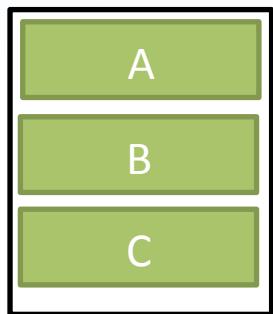
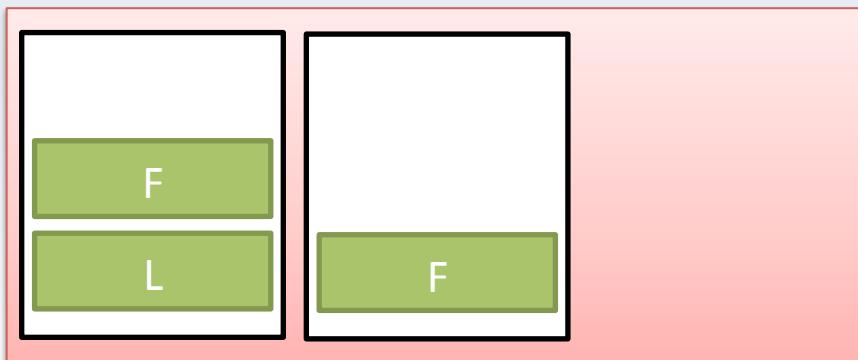
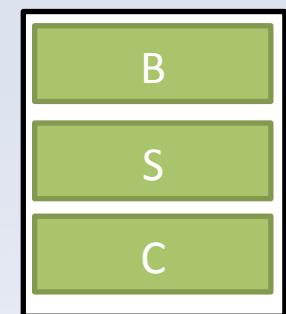
block 3



block 4

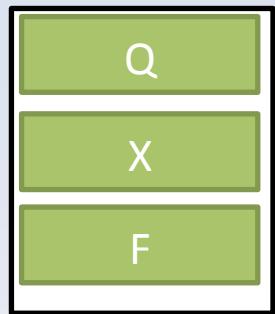


block 5

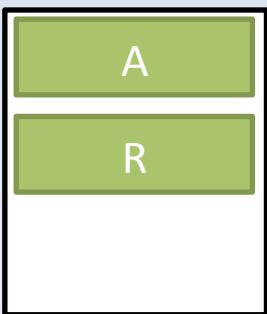


# Phase 2

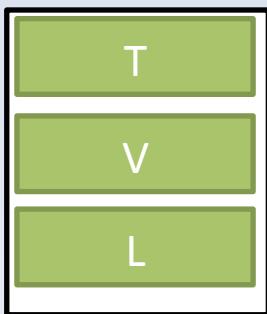
block 0



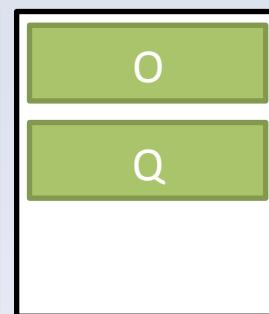
block 1



block 2



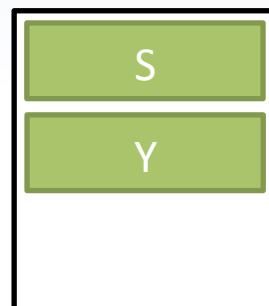
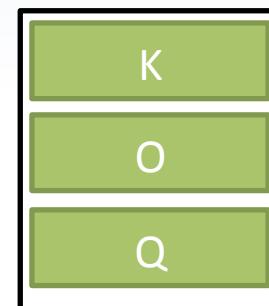
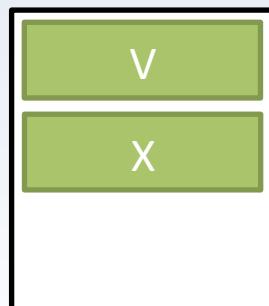
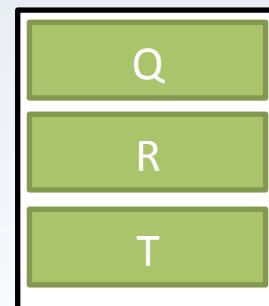
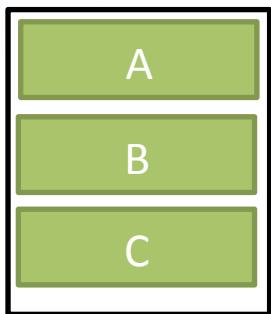
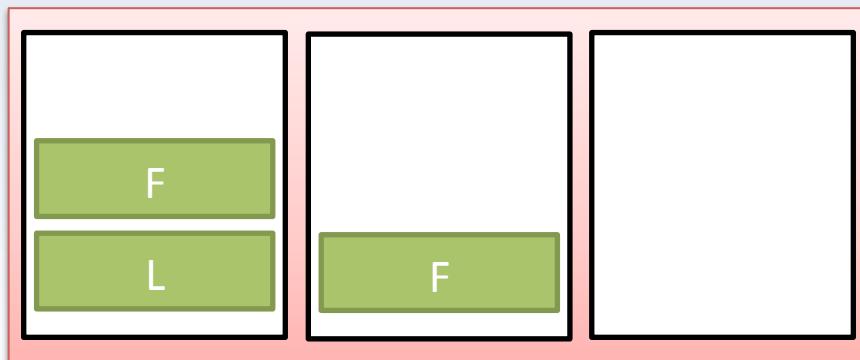
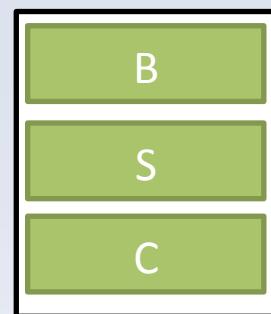
block 3



block 4

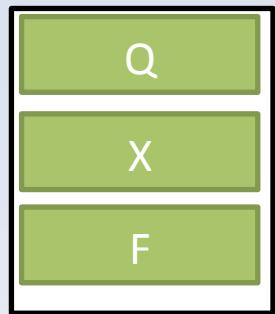


block 5

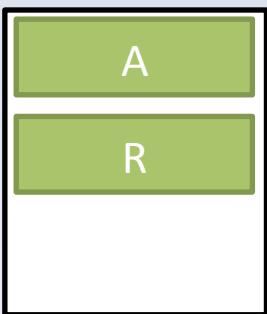


# Phase 2

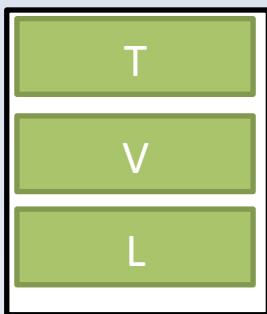
block 0



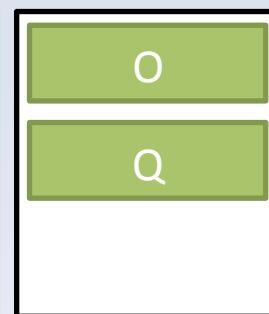
block 1



block 2



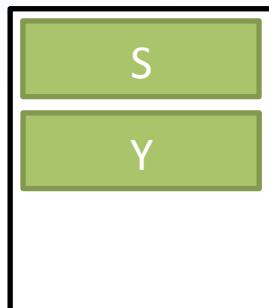
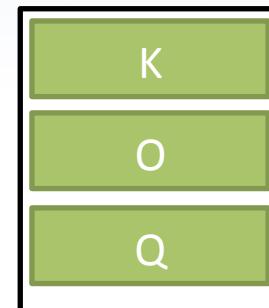
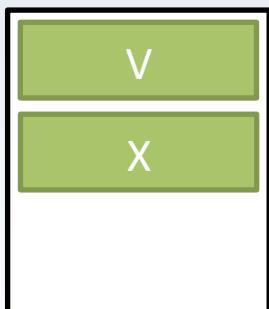
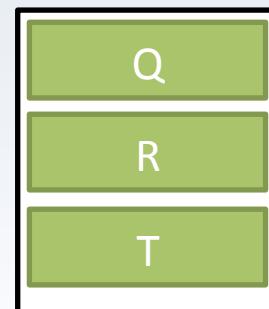
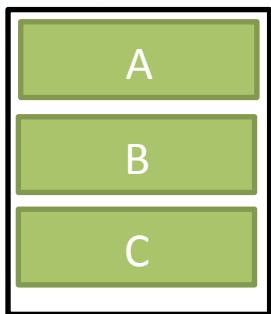
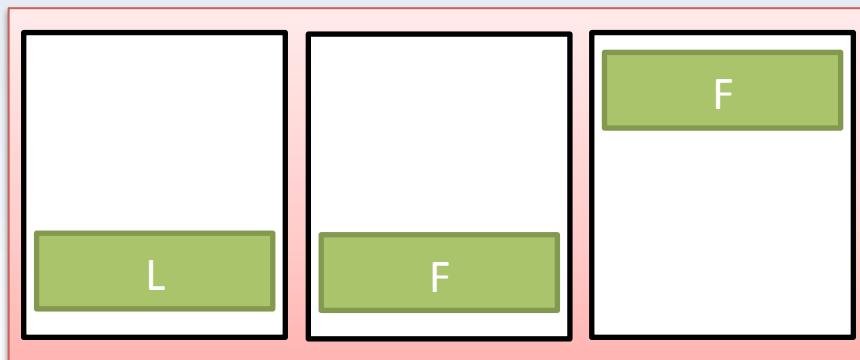
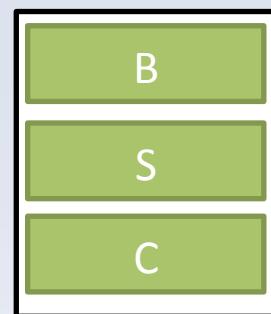
block 3



block 4

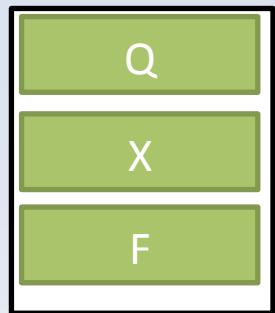


block 5

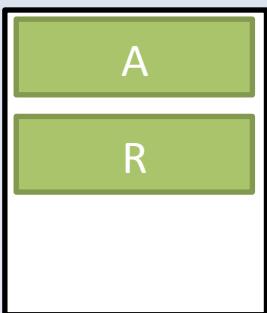


# Phase 2

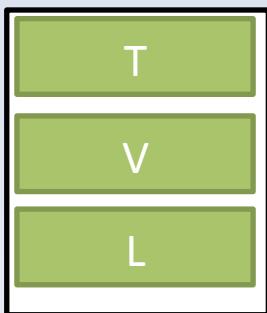
block 0



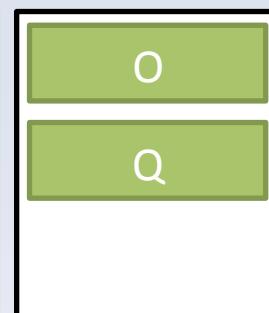
block 1



block 2



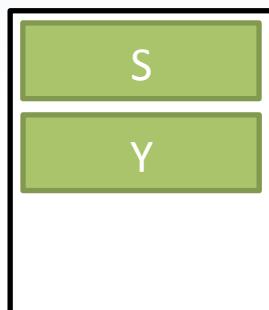
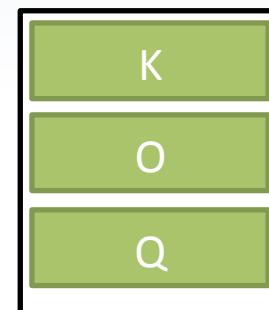
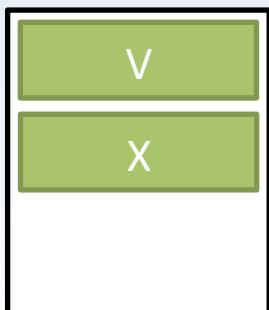
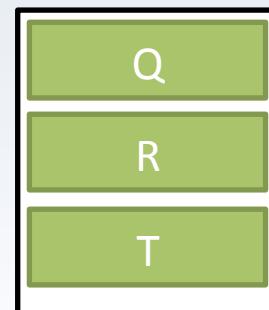
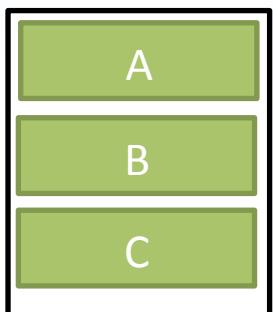
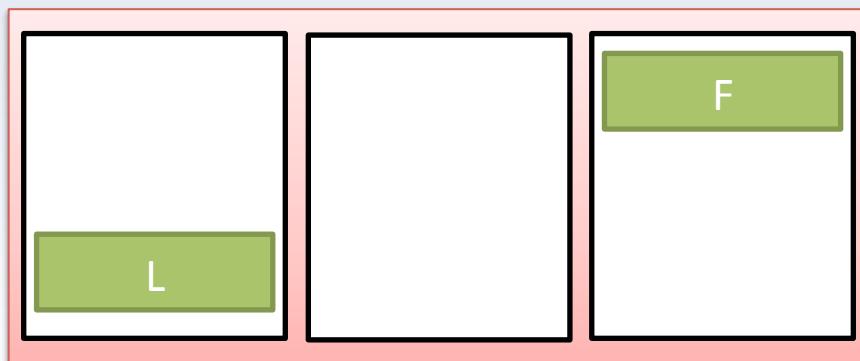
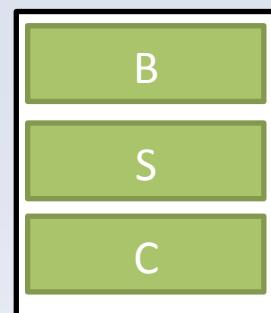
block 3



block 4

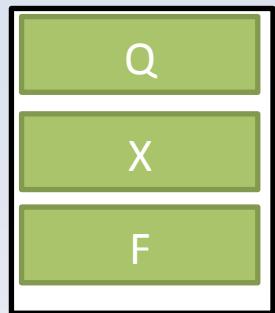


block 5

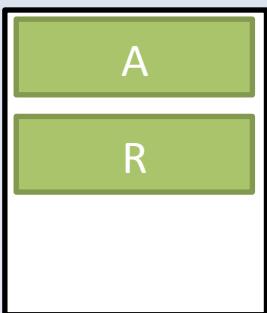


# Phase 2

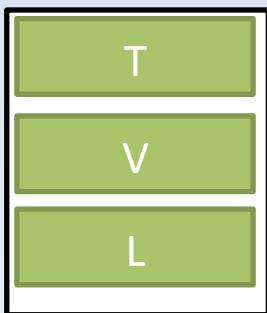
block 0



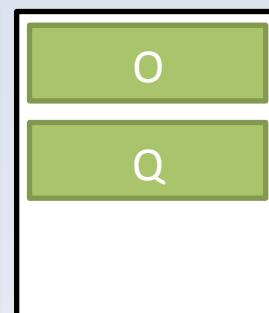
block 1



block 2



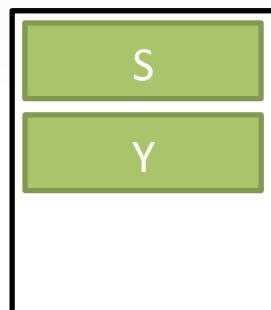
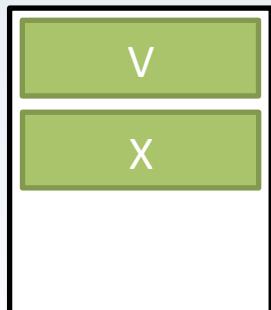
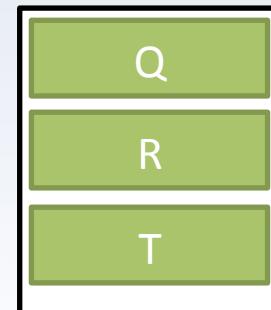
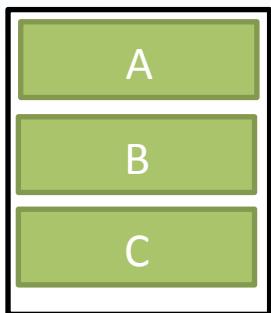
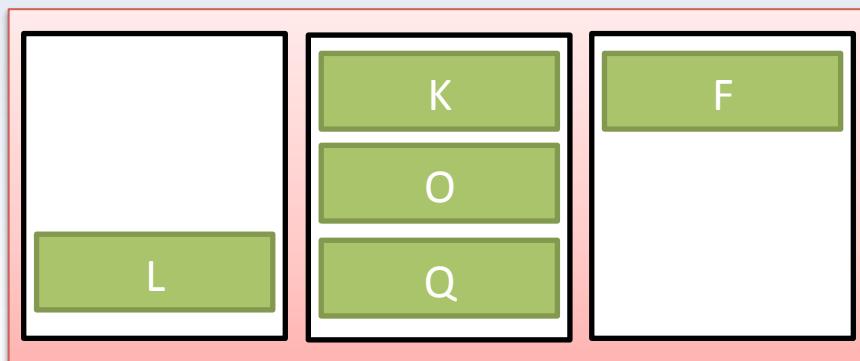
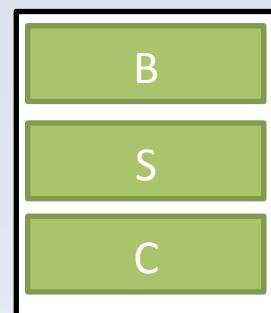
block 3



block 4

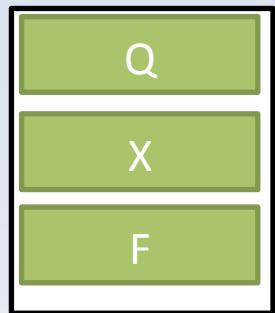


block 5

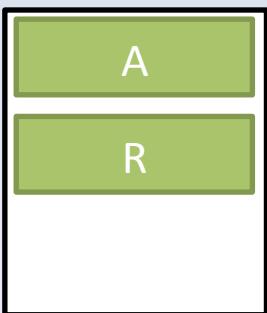


# Phase 2

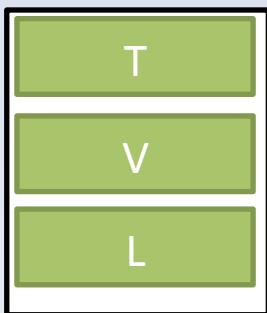
block 0



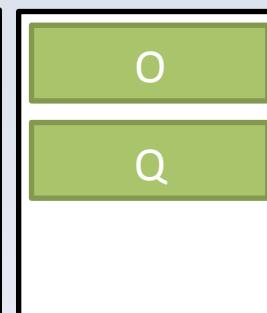
block 1



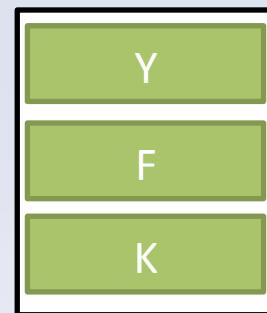
block 2



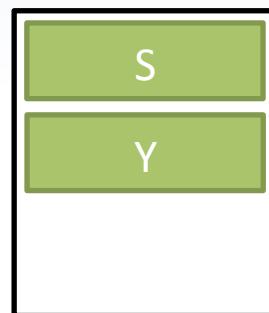
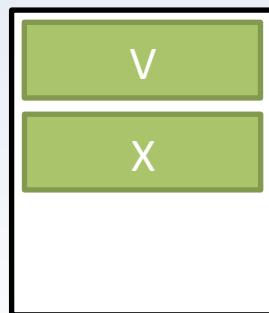
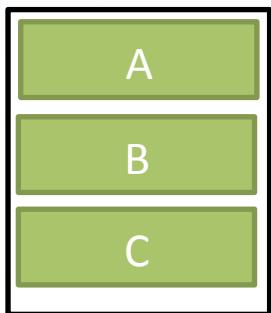
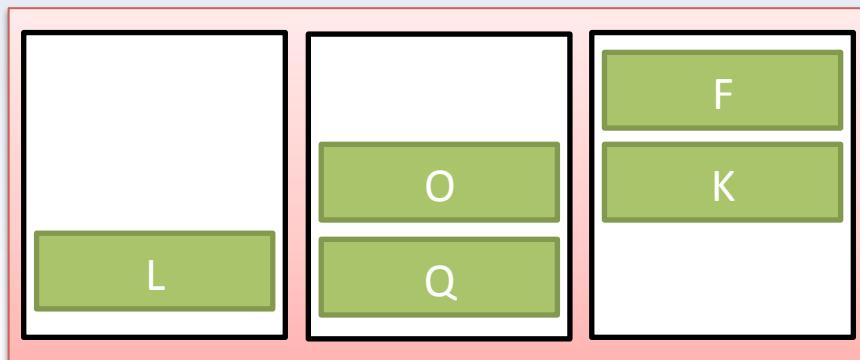
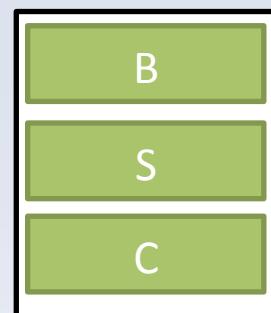
block 3



block 4

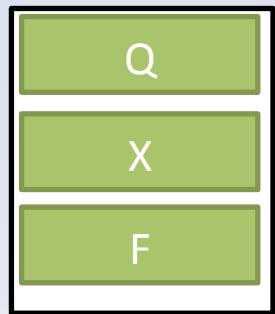


block 5

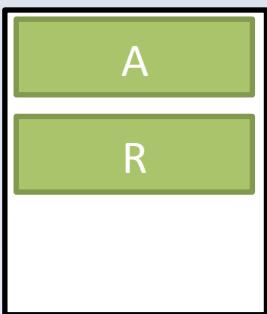


# Phase 2

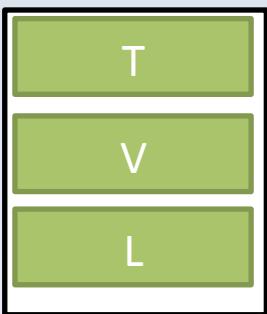
block 0



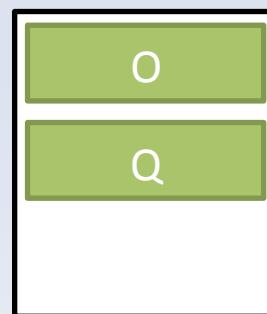
block 1



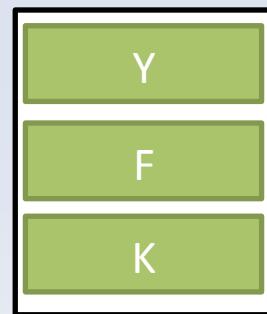
block 2



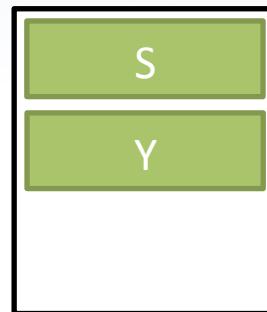
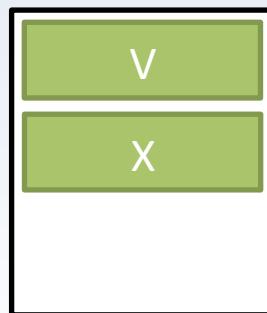
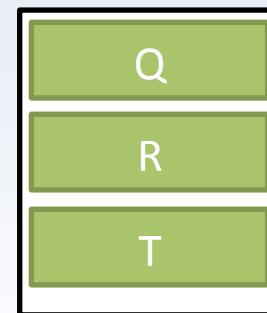
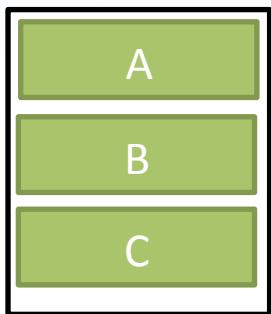
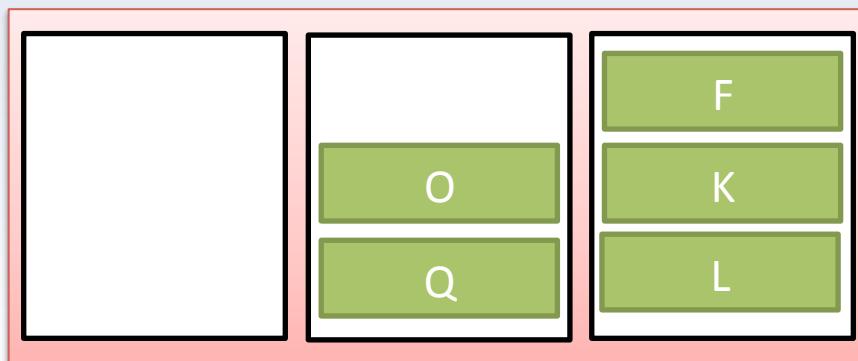
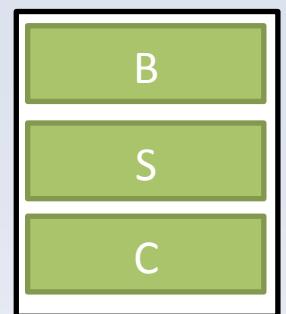
block 3



block 4

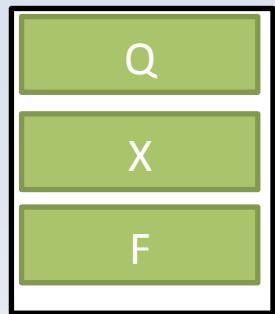


block 5

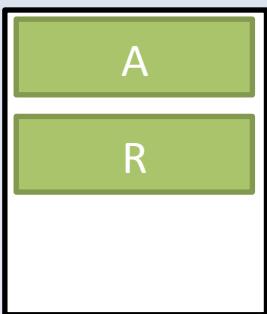


# Phase 2

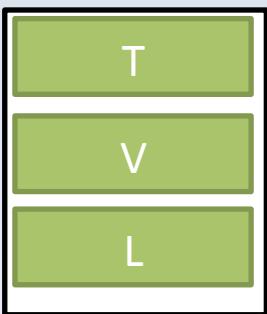
block 0



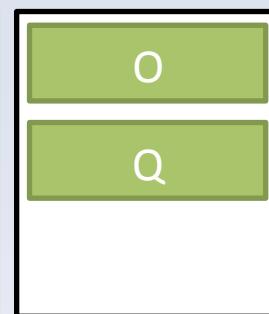
block 1



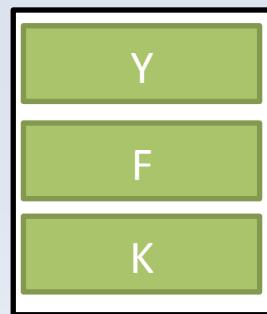
block 2



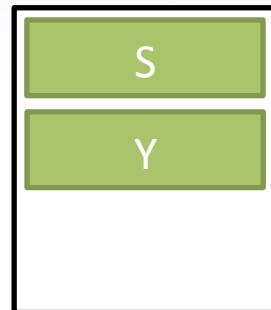
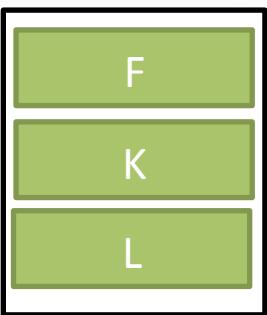
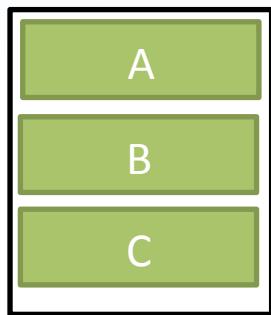
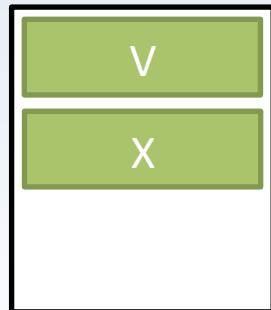
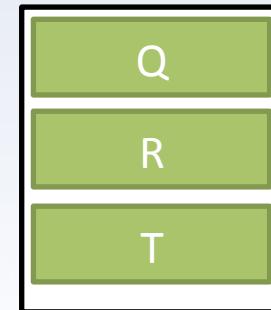
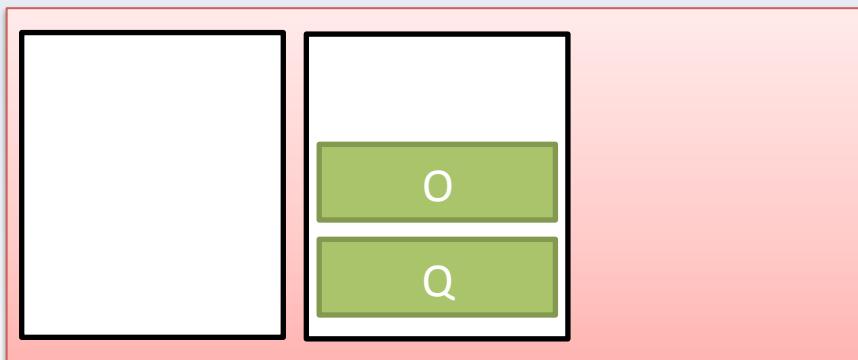
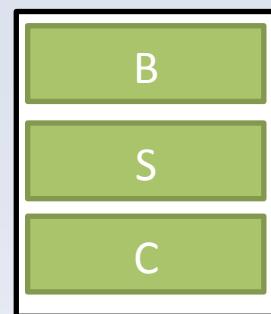
block 3



block 4

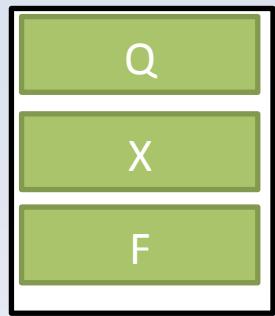


block 5

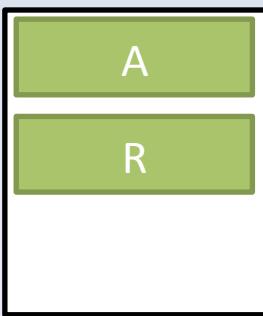


# Phase 2

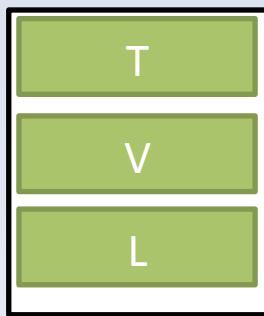
block 0



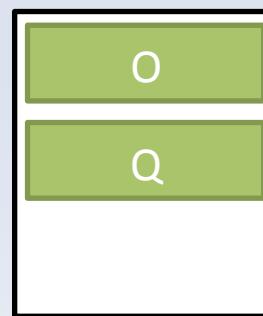
block 1



block 2



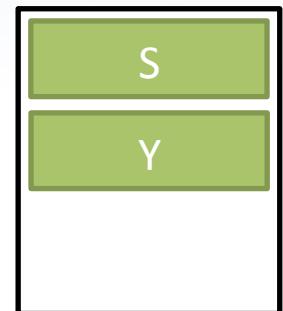
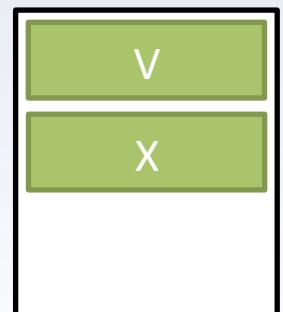
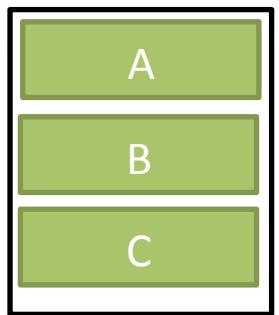
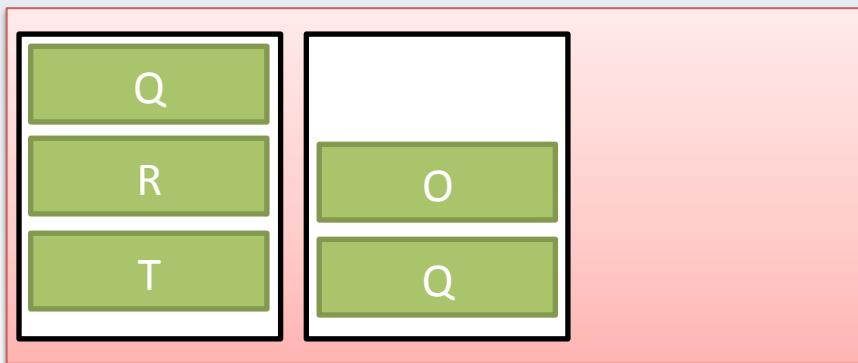
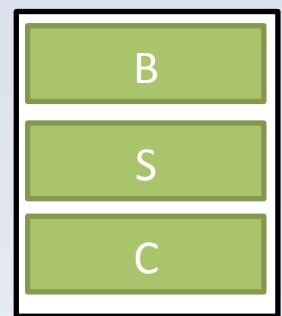
block 3



block 4

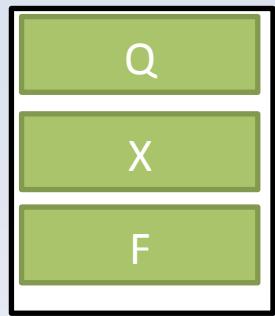


block 5

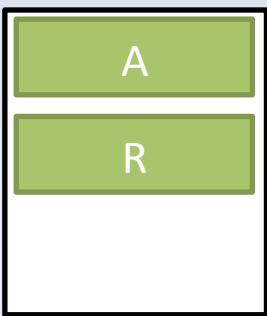


# Phase 2

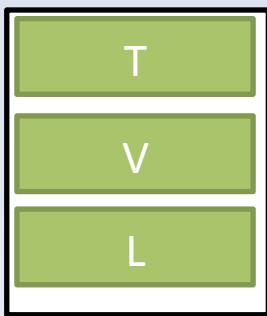
block 0



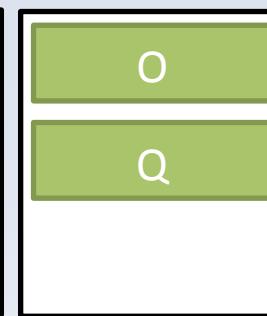
block 1



block 2



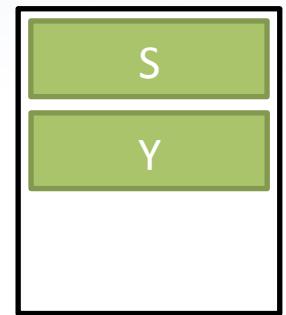
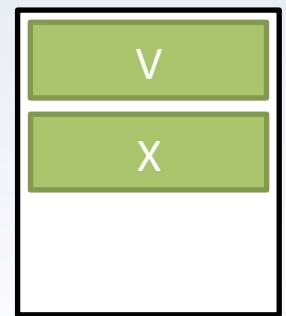
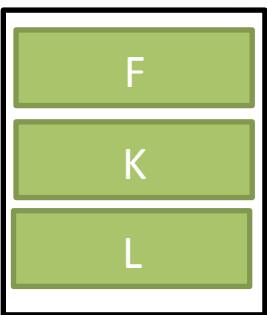
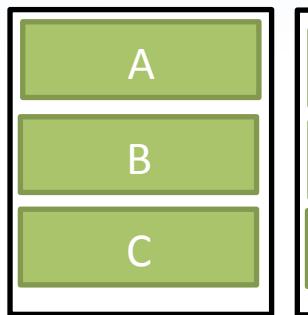
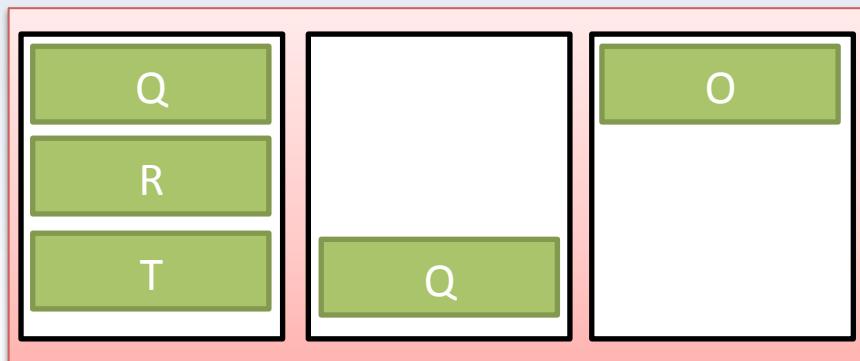
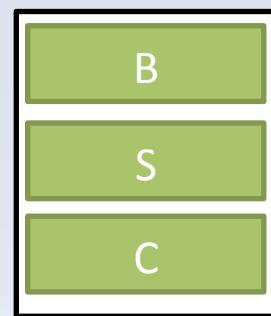
block 3



block 4

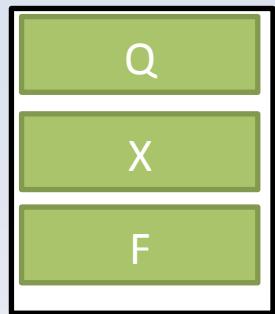


block 5

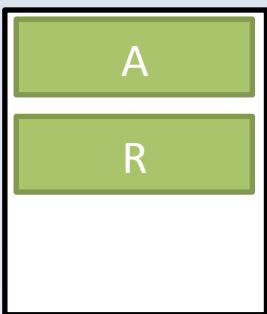


# Phase 2

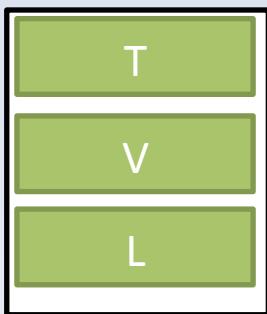
block 0



block 1



block 2



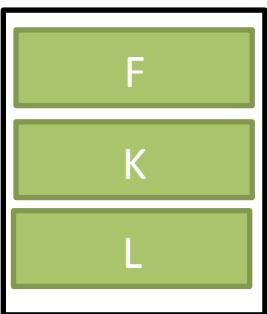
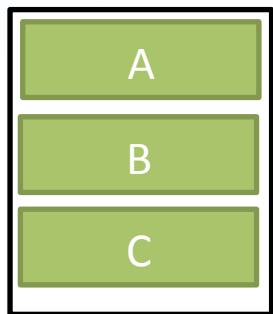
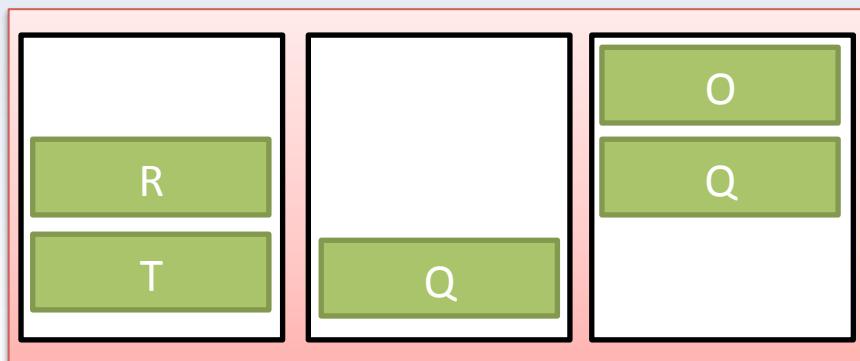
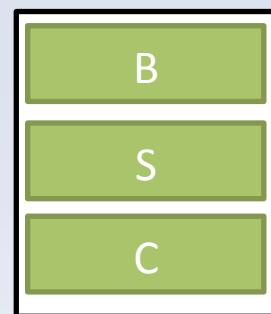
block 3



block 4

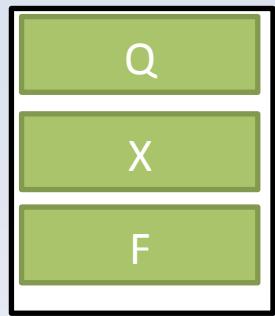


block 5

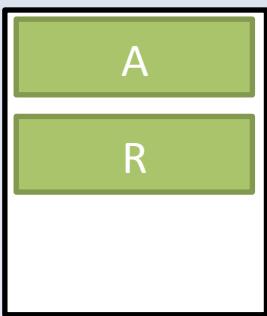


# Phase 2

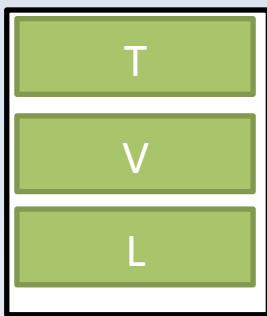
block 0



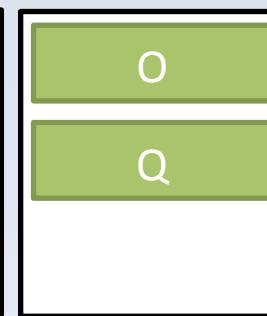
block 1



block 2



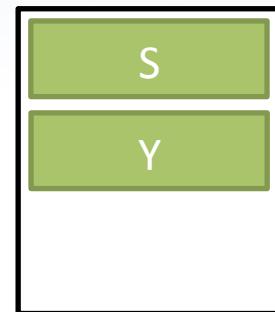
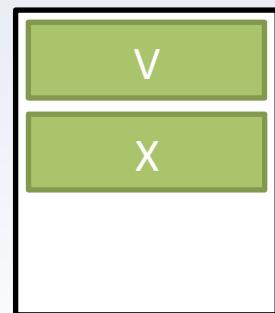
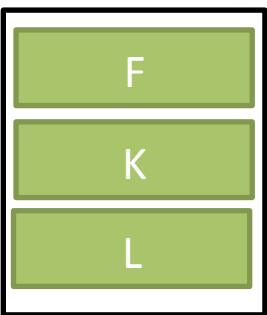
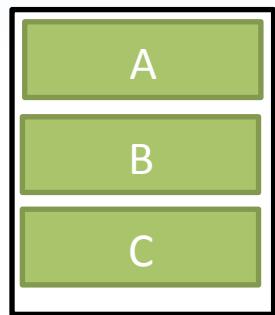
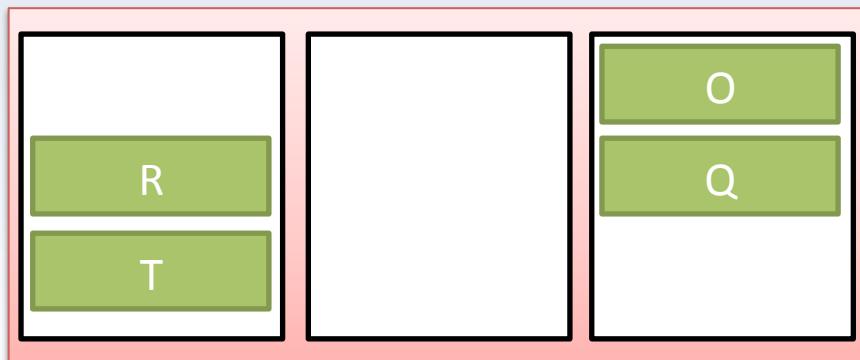
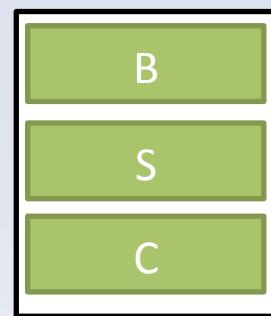
block 3



block 4

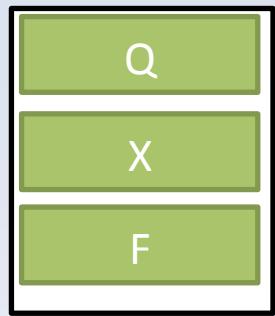


block 5

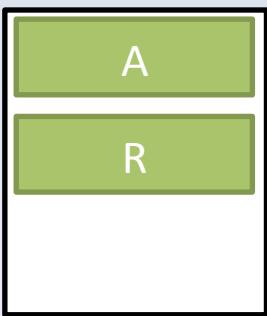


# Phase 2

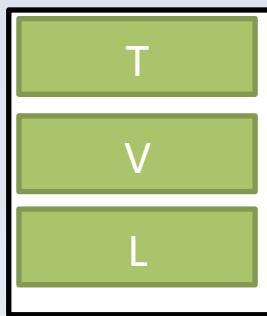
block 0



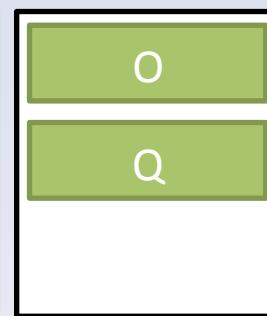
block 1



block 2



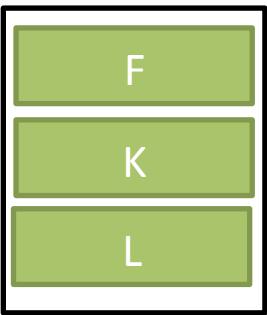
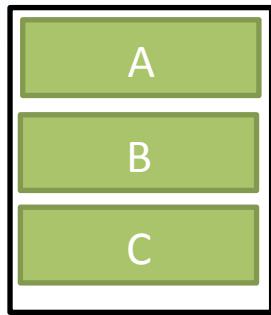
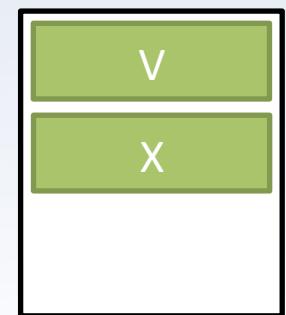
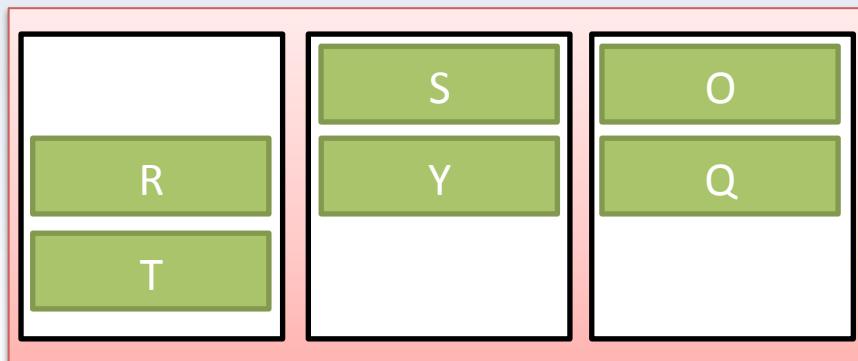
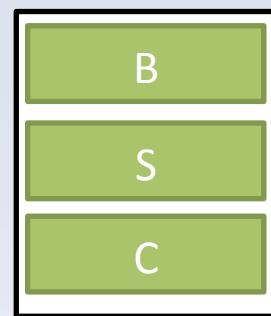
block 3



block 4

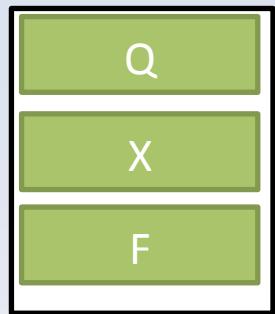


block 5

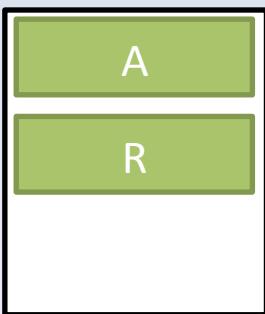


# Phase 2

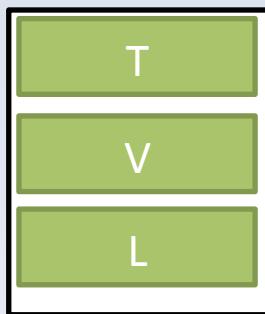
block 0



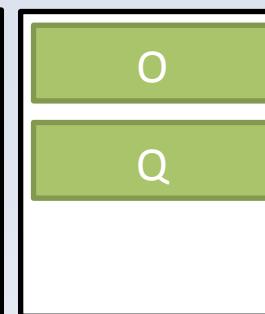
block 1



block 2



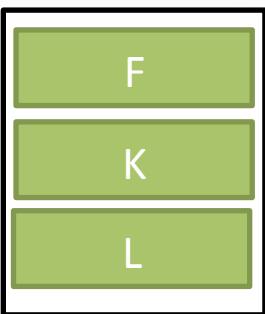
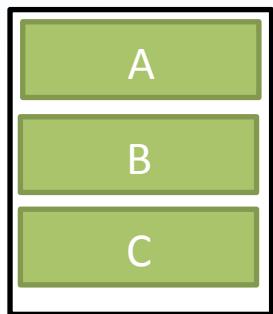
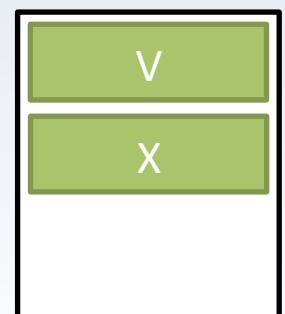
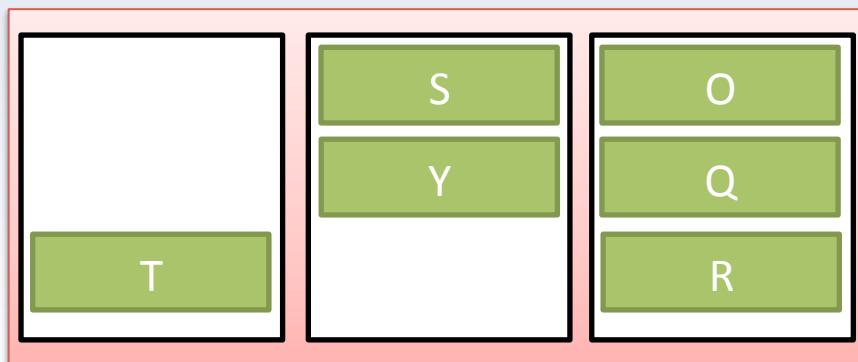
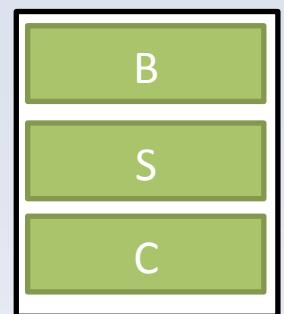
block 3



block 4

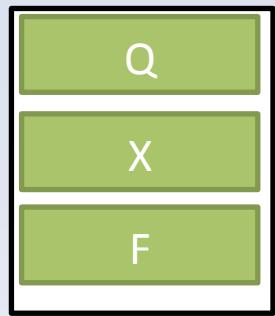


block 5

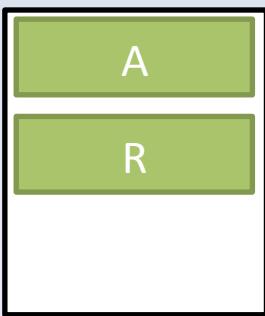


# Phase 2

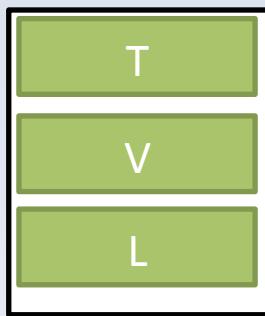
block 0



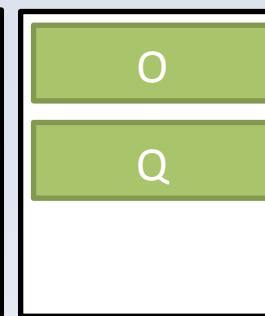
block 1



block 2



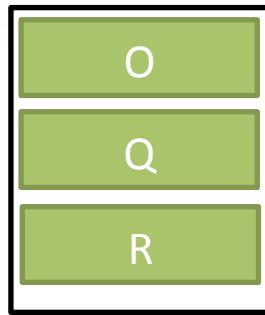
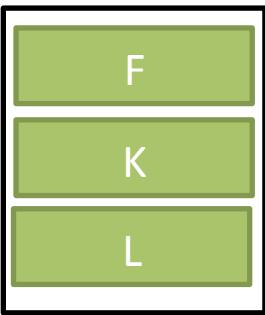
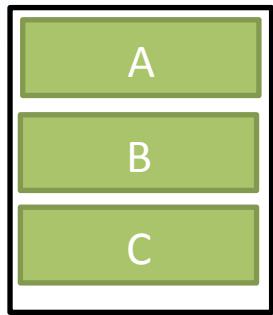
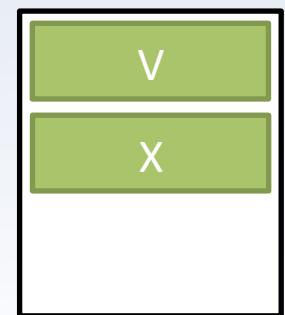
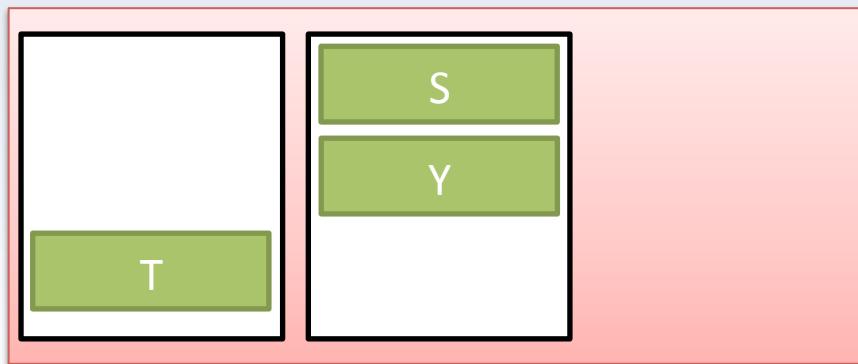
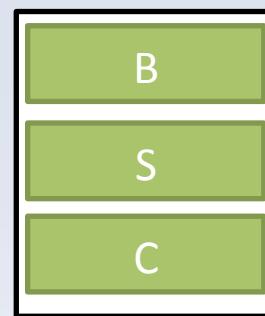
block 3



block 4

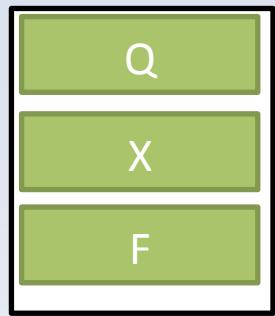


block 5

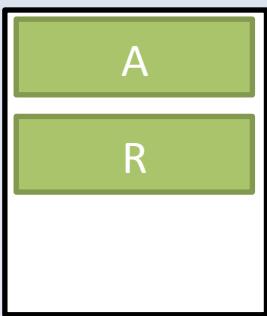


# Phase 2

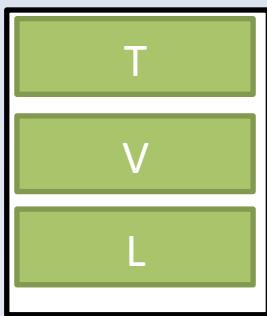
block 0



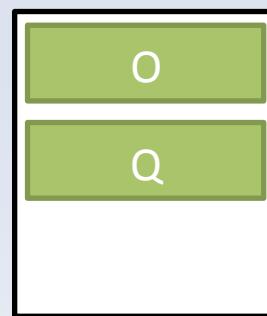
block 1



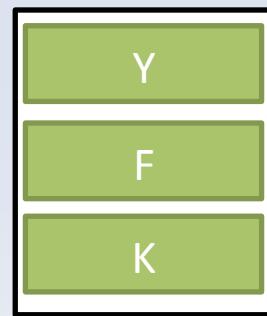
block 2



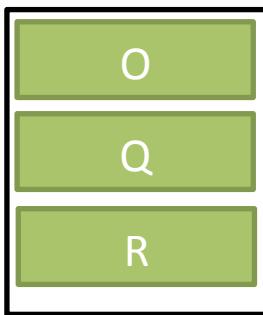
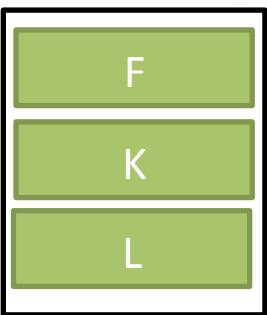
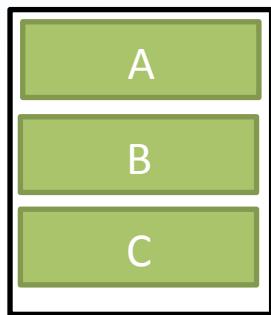
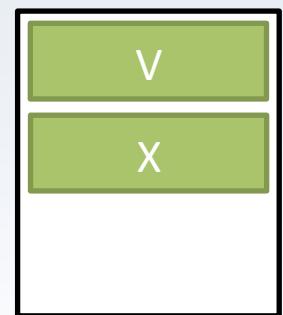
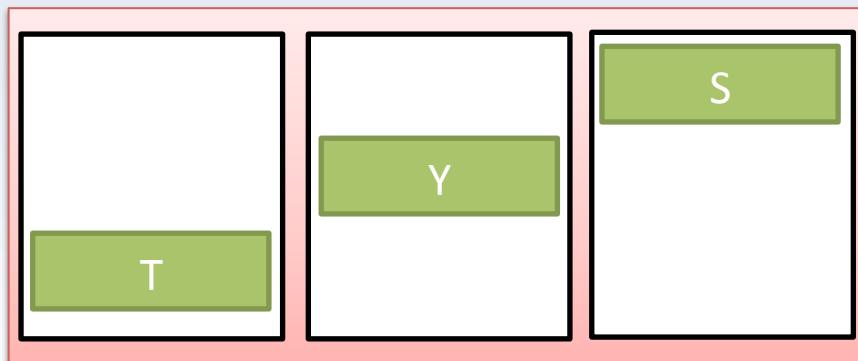
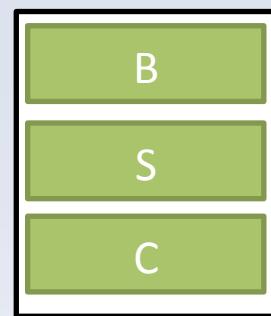
block 3



block 4

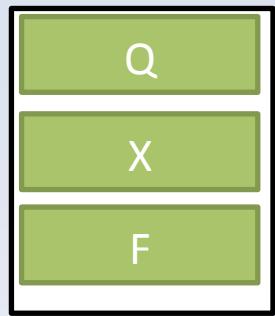


block 5

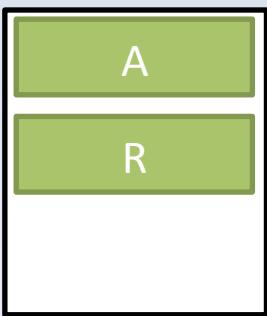


# Phase 2

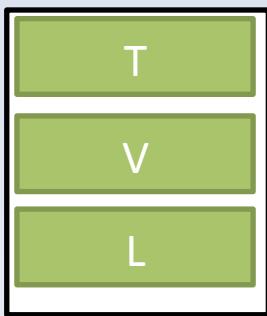
block 0



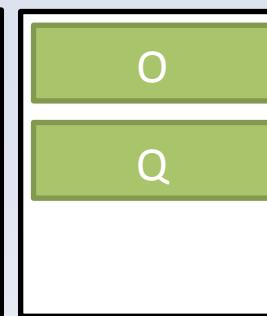
block 1



block 2



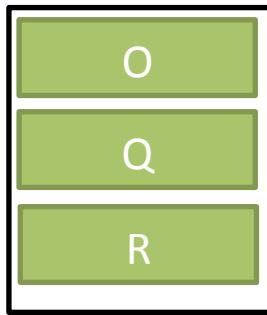
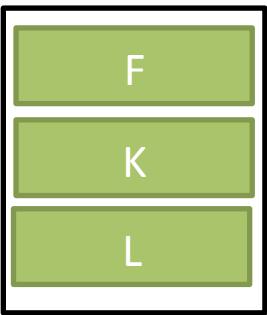
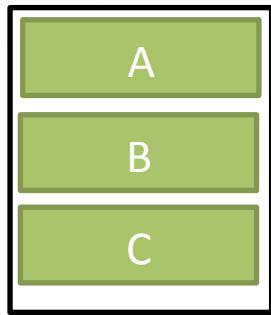
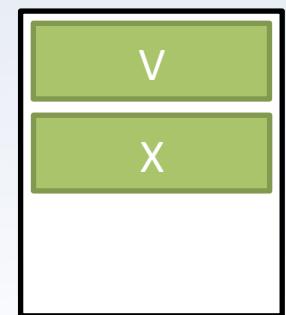
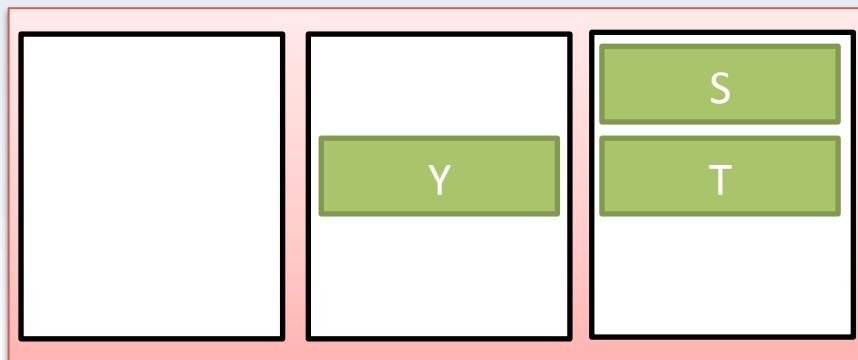
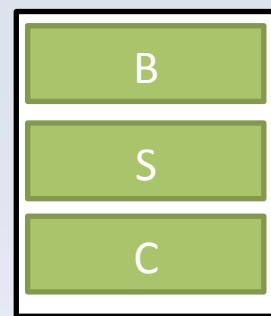
block 3



block 4

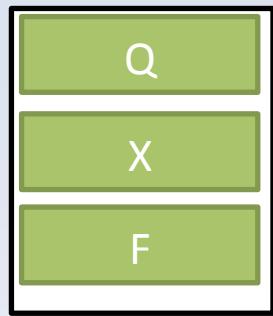


block 5

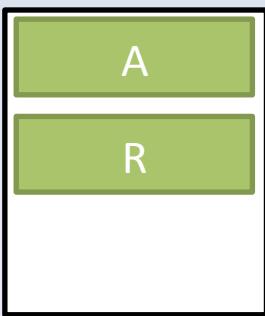


# Phase 2

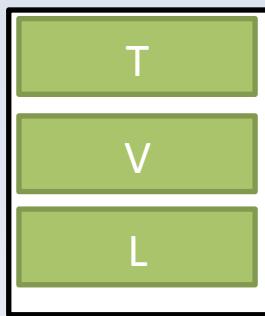
block 0



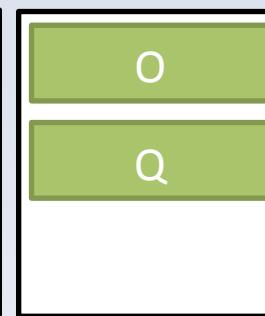
block 1



block 2



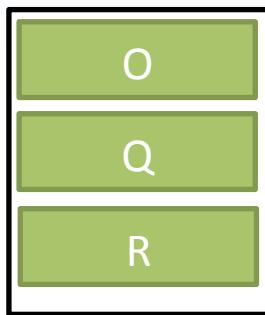
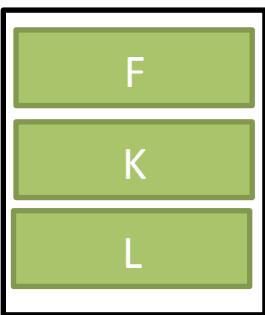
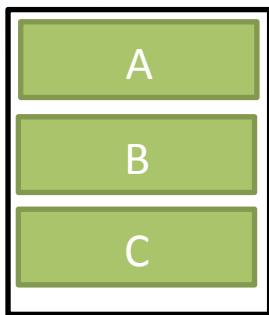
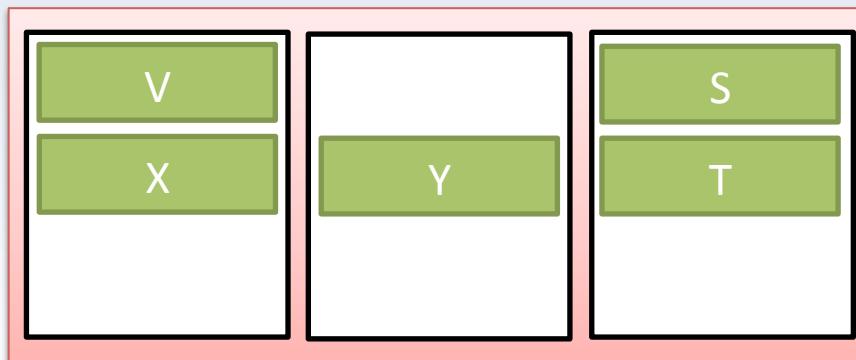
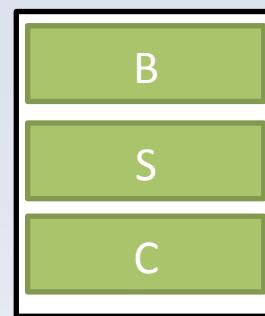
block 3



block 4

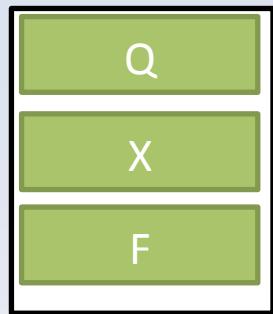


block 5

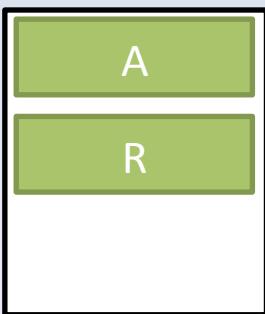


# Phase 2

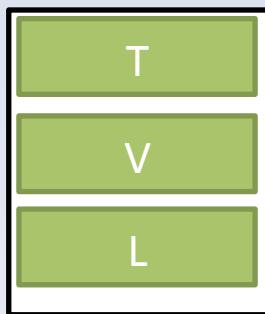
block 0



block 1



block 2



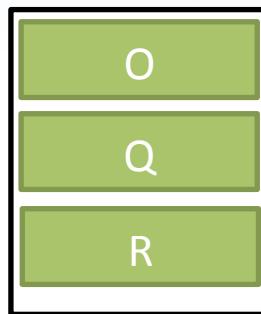
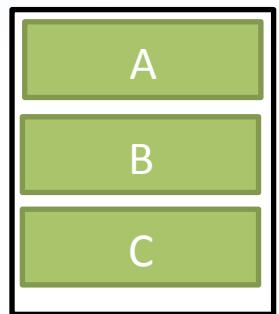
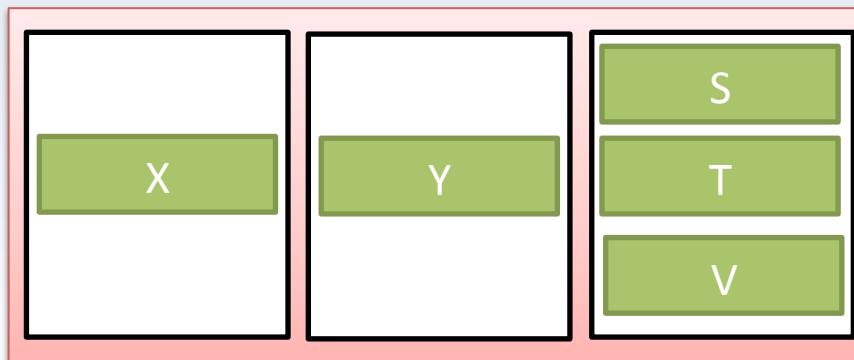
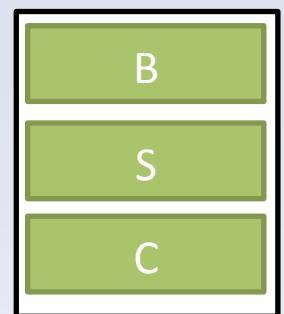
block 3



block 4

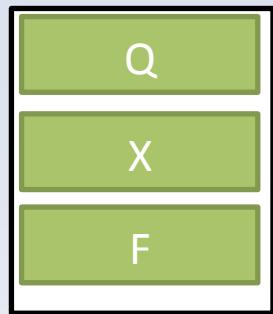


block 5

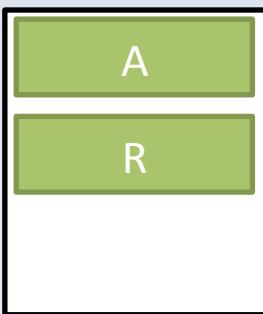


# Phase 2

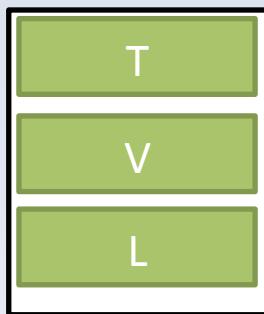
block 0



block 1



block 2



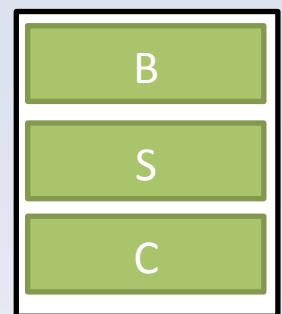
block 3



block 4



block 5



Q

X

F

A

R

T

V

L

O

Q

Y

F

K

B

S

C

X

Y

A

B

C

F

K

L

O

Q

R

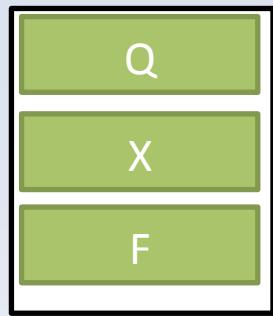
S

T

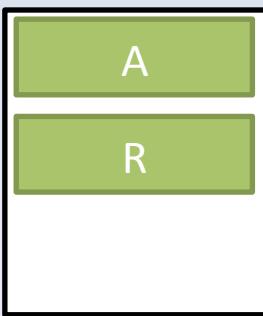
V

# Phase 2

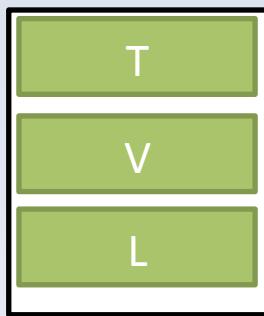
block 0



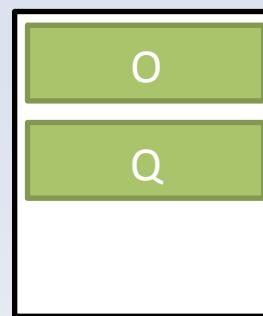
block 1



block 2



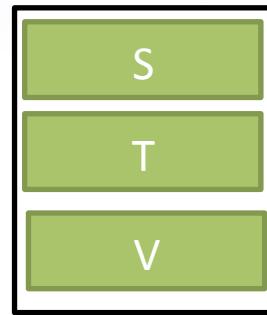
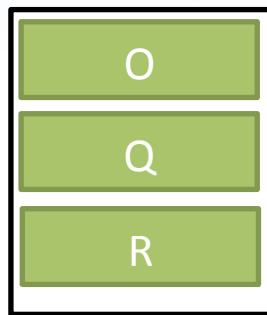
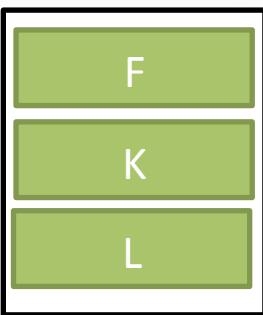
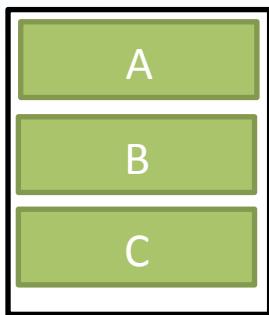
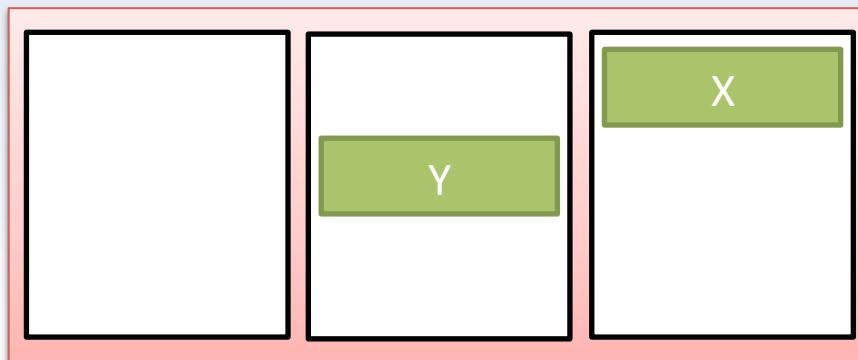
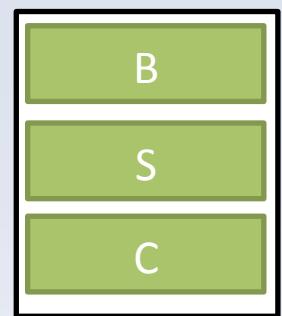
block 3



block 4

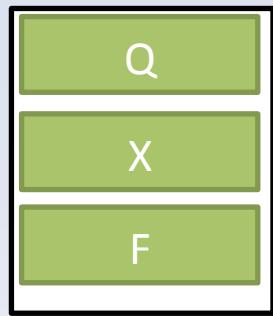


block 5

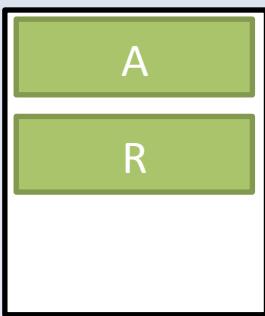


# Phase 2

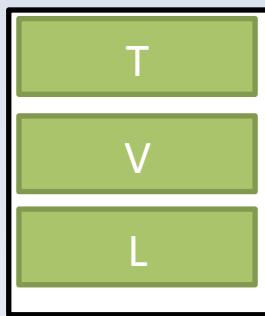
block 0



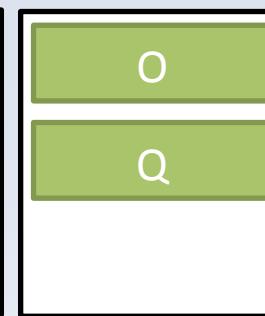
block 1



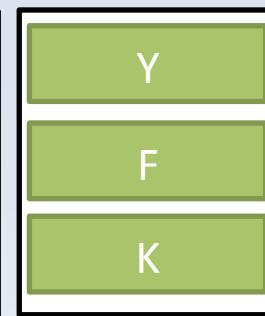
block 2



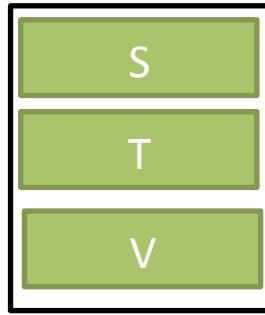
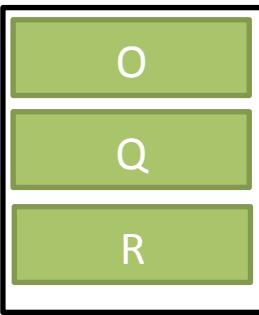
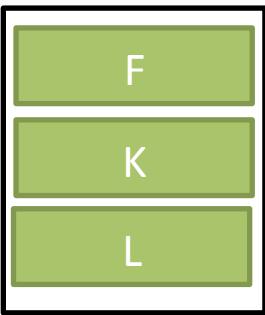
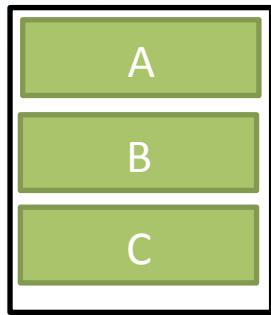
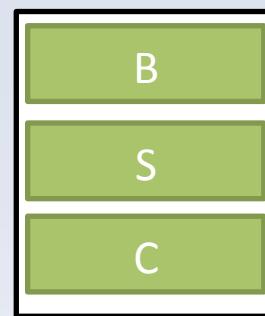
block 3



block 4

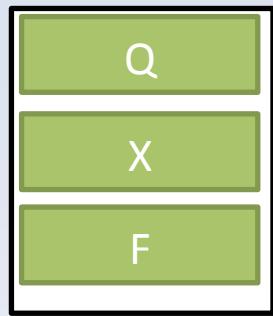


block 5

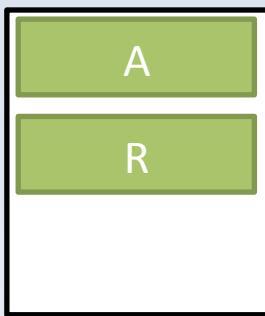


# Phase 2

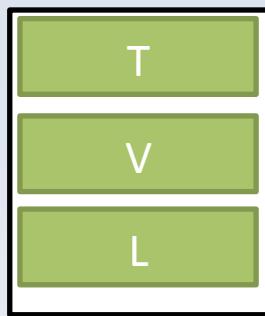
block 0



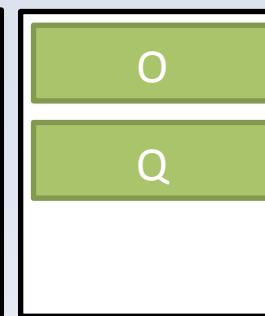
block 1



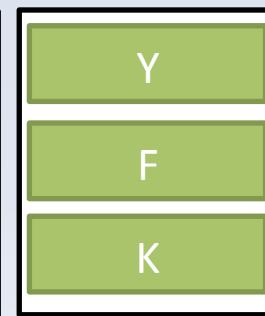
block 2



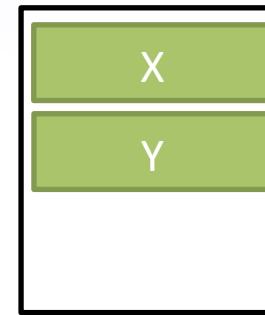
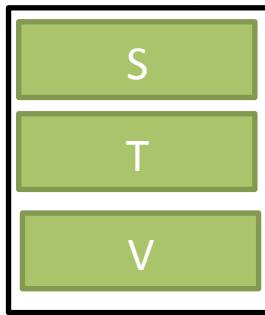
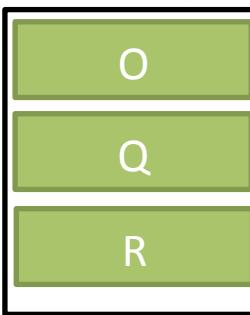
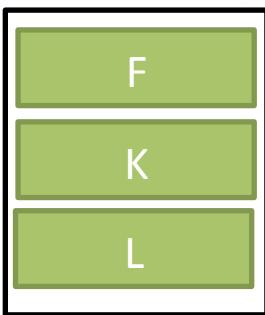
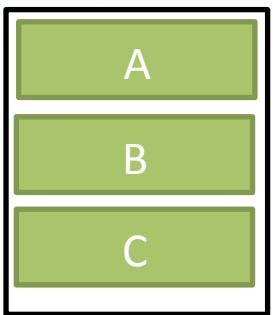
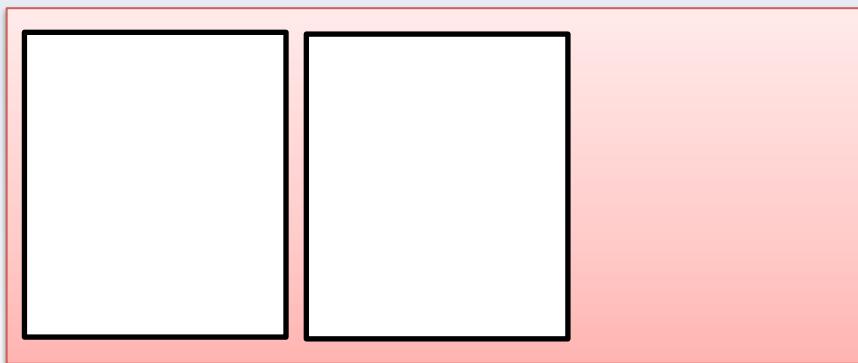
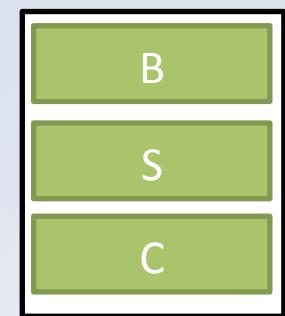
block 3

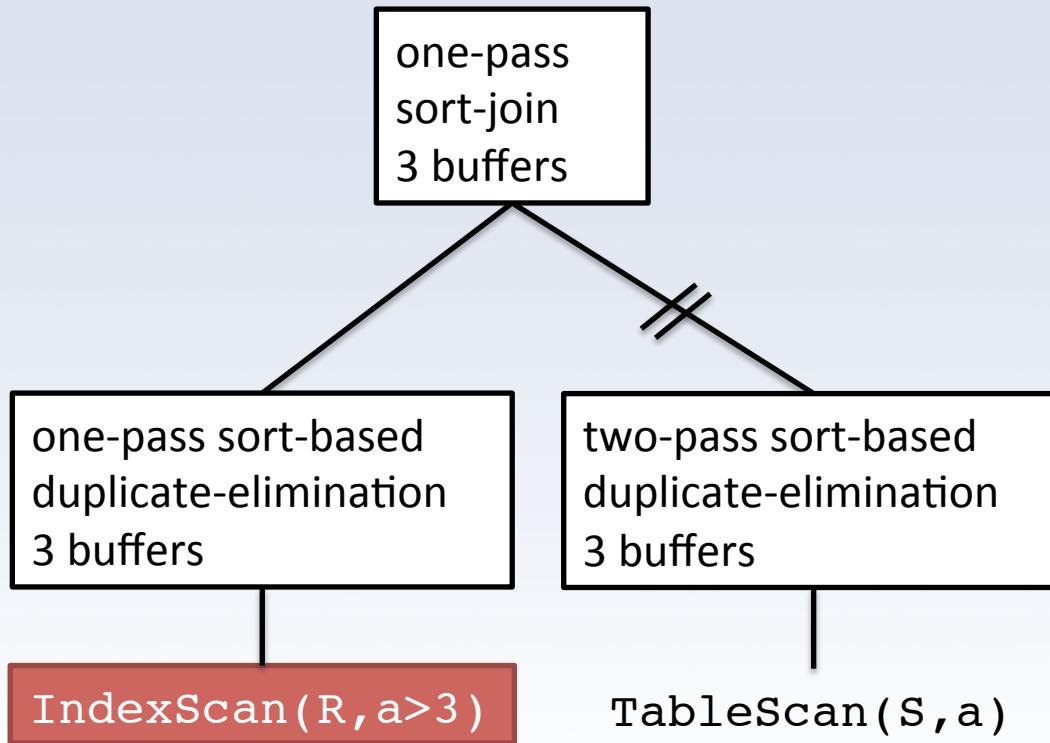


block 4

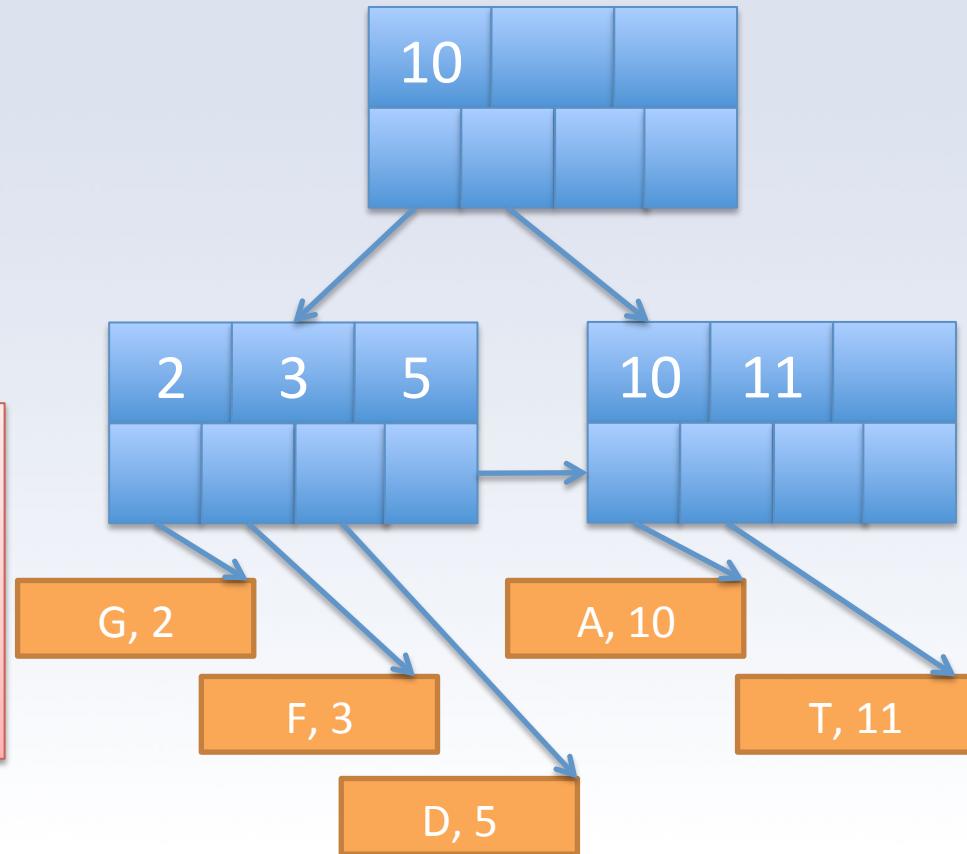
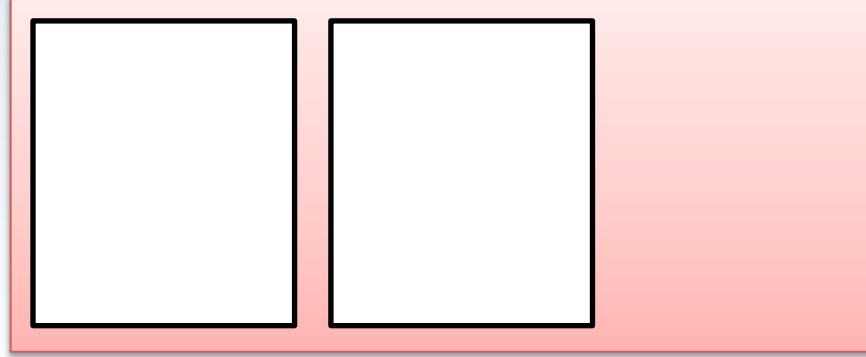


block 5

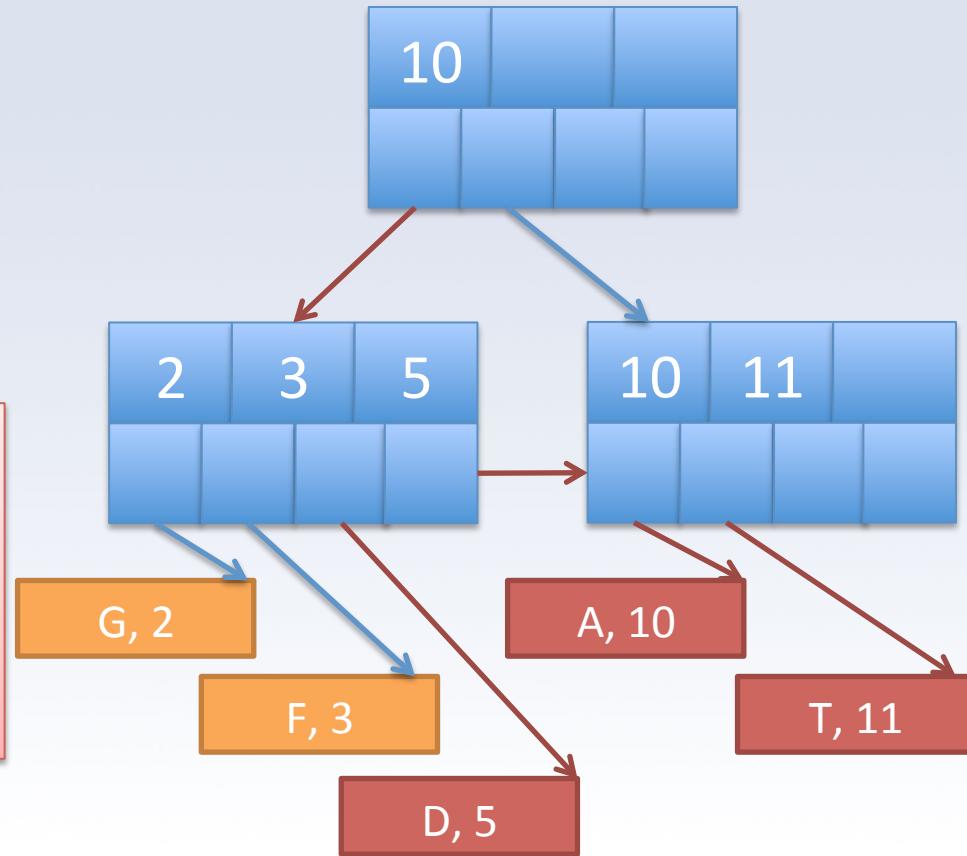
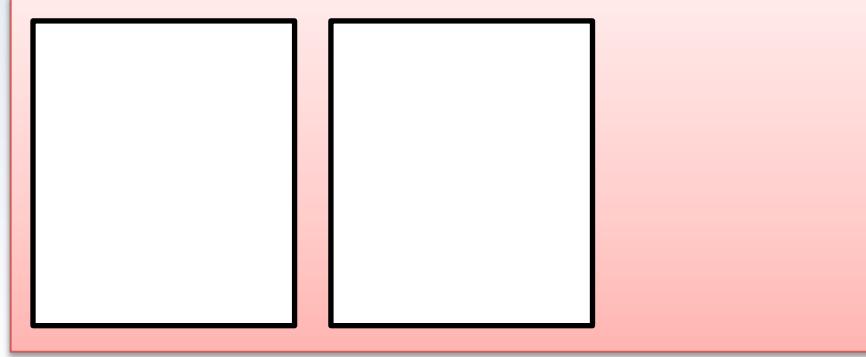




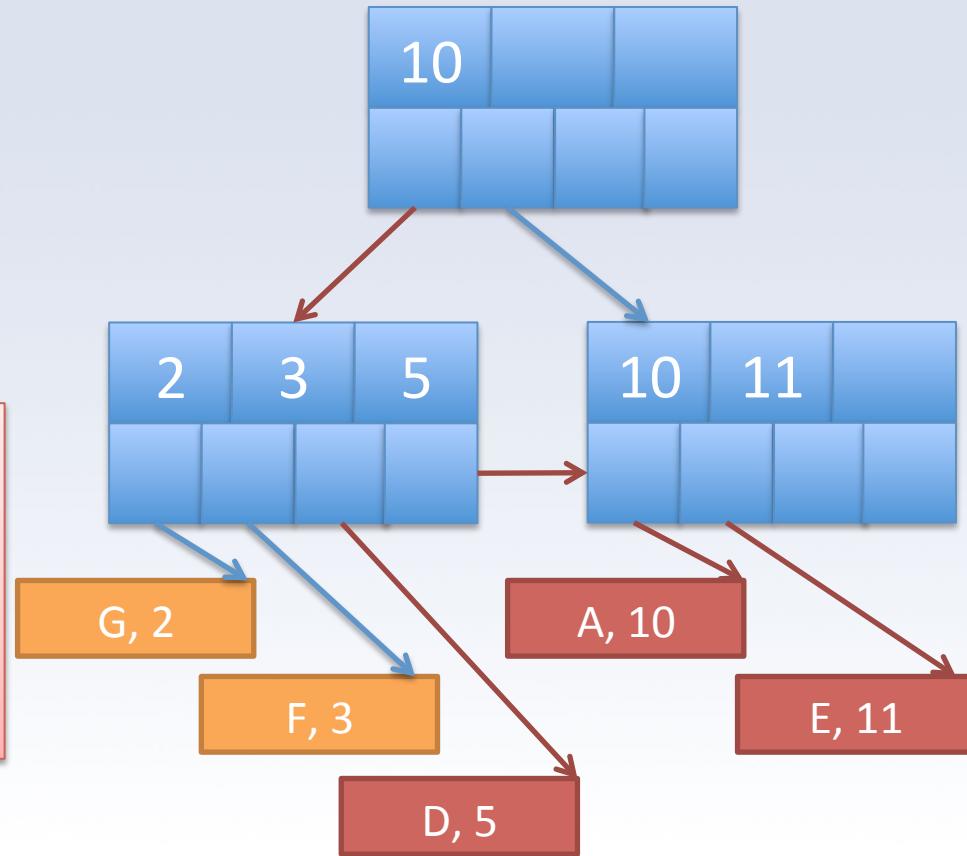
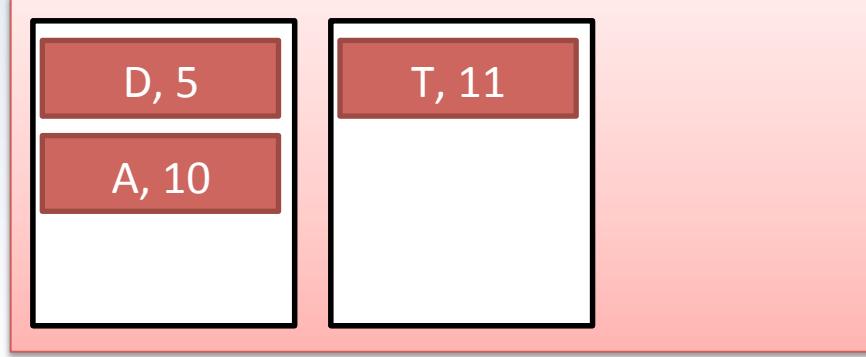
# Index Scan

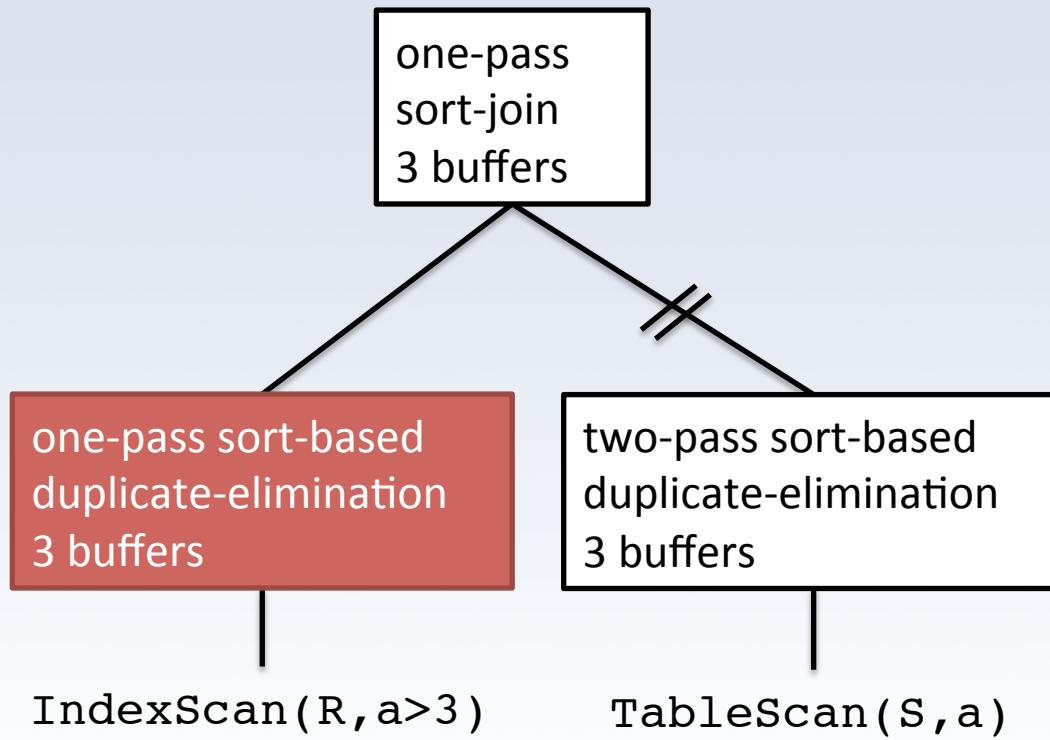


# Index Scan

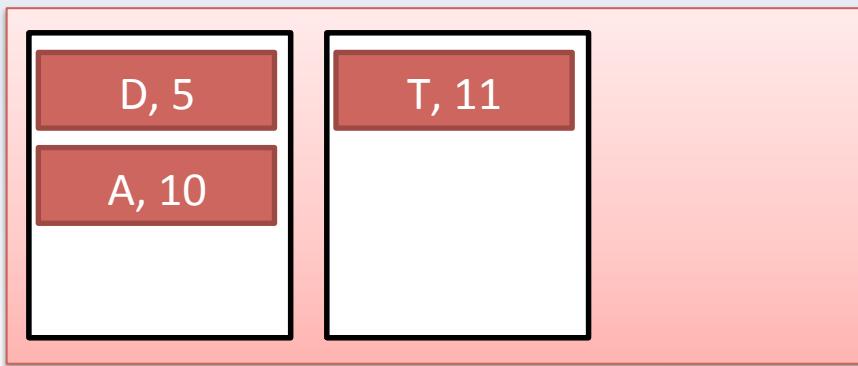


# Index Scan





# One Pass Sort



# One Pass Sort

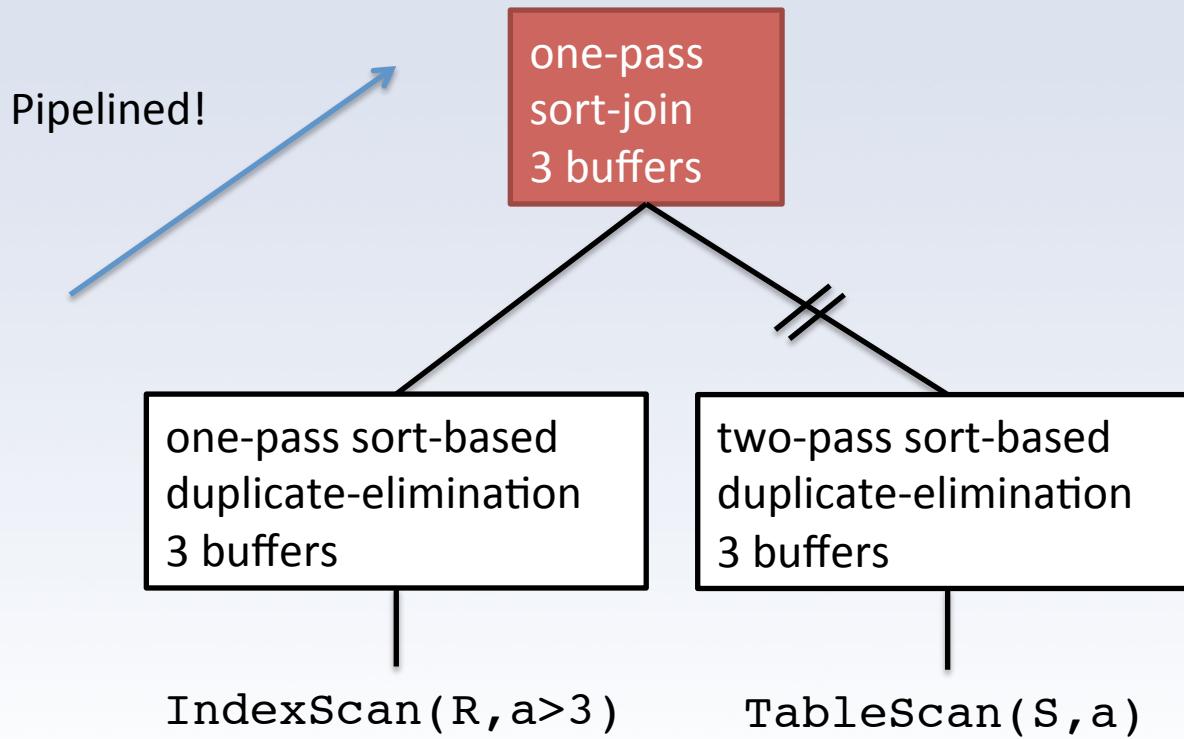
A, 10

D, 5

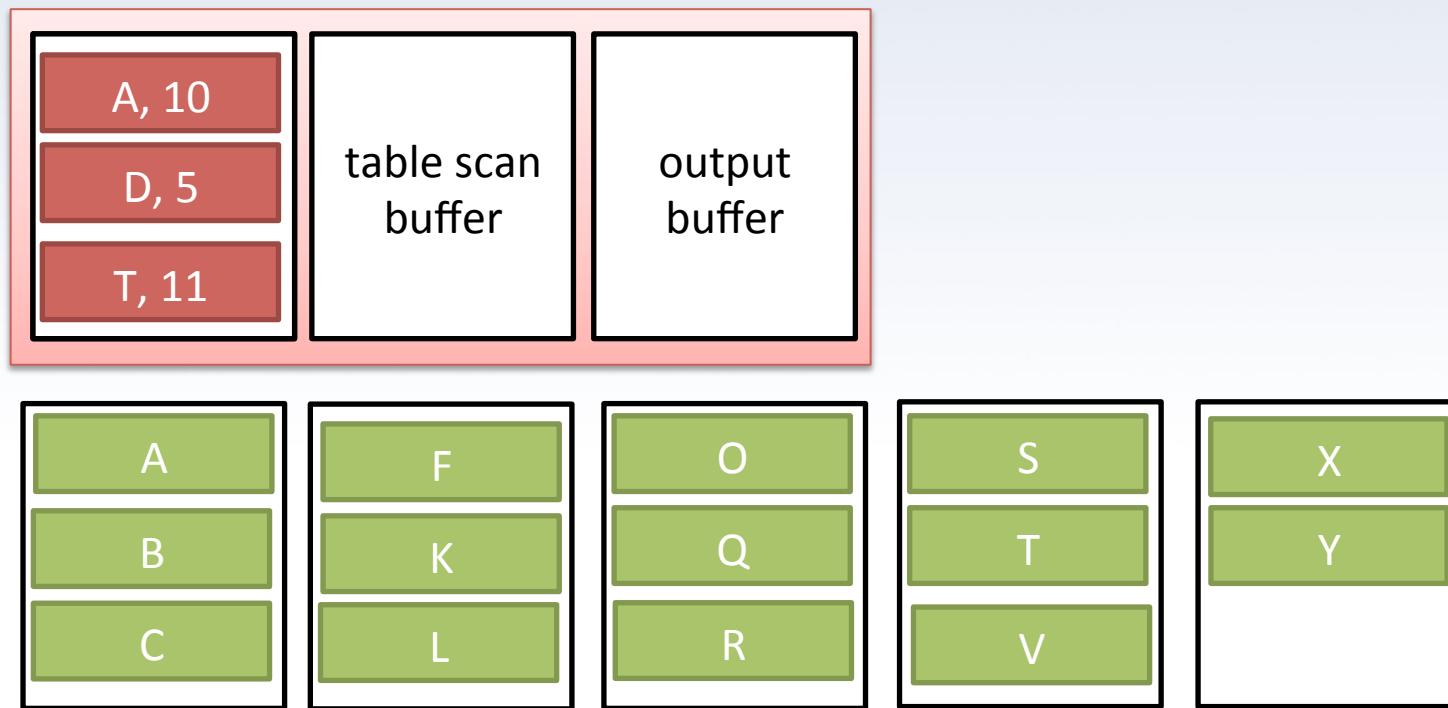
T, 11



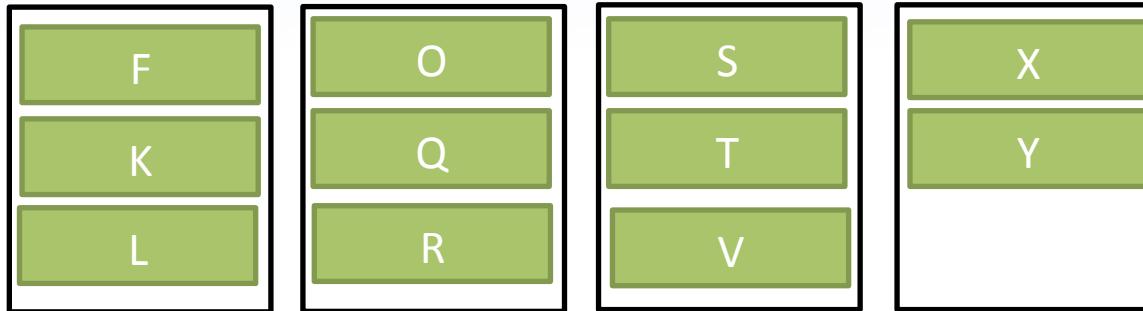
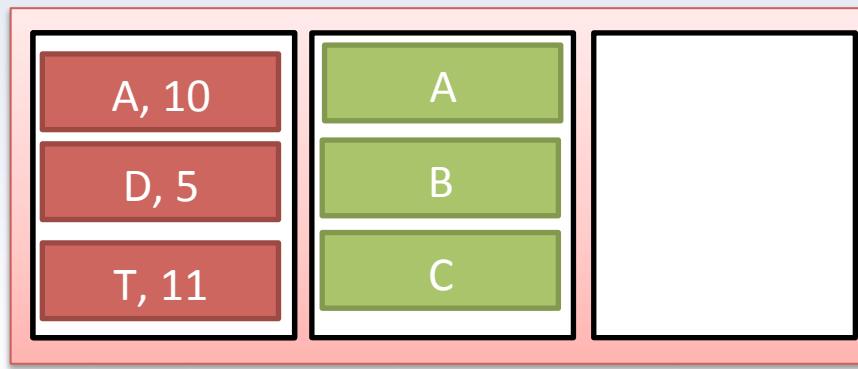
# Physical Query Plan



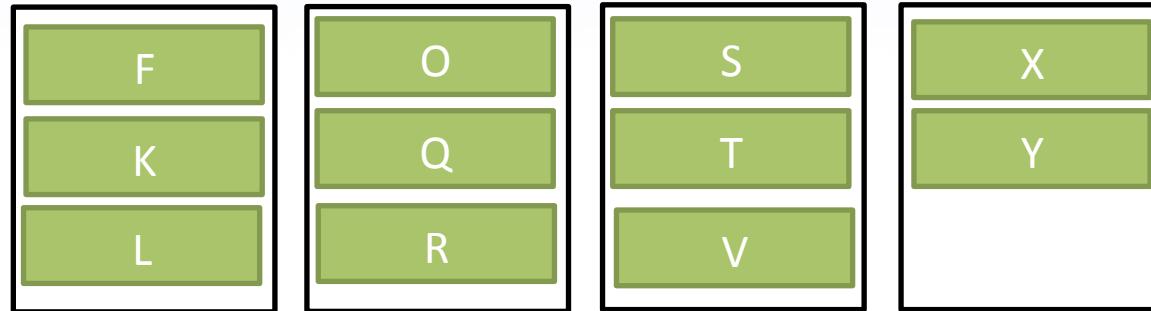
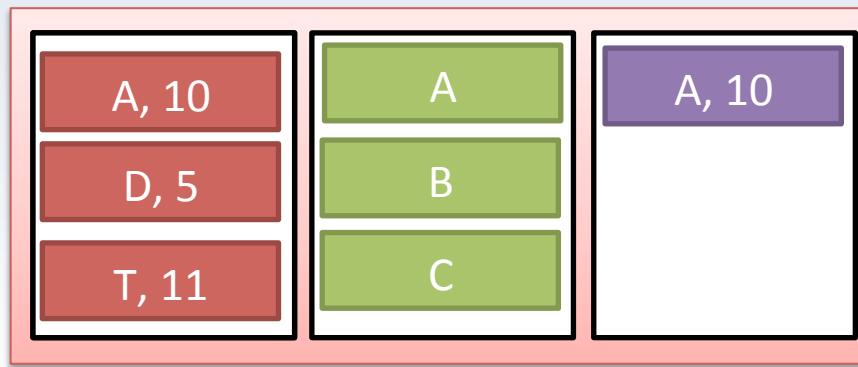
# One Pass Sort-Join



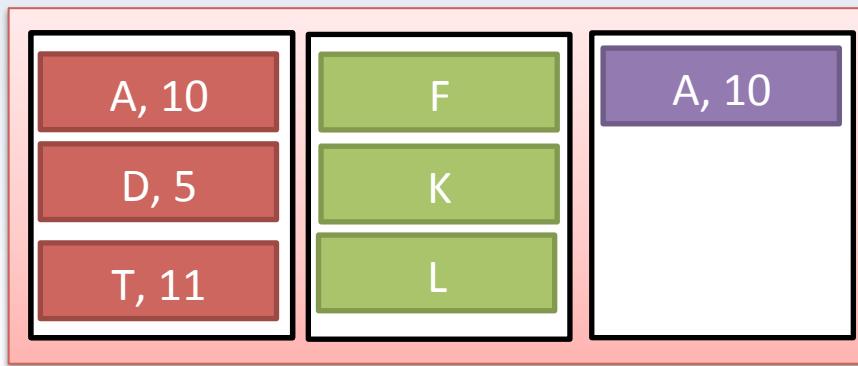
# One Pass Sort-Join



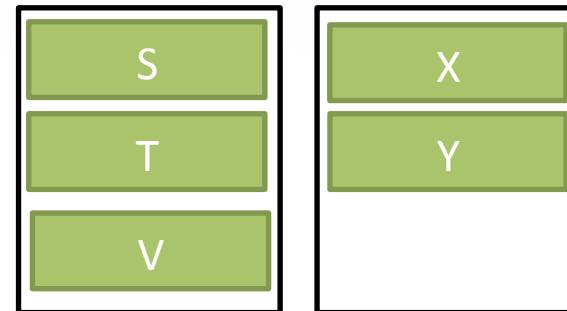
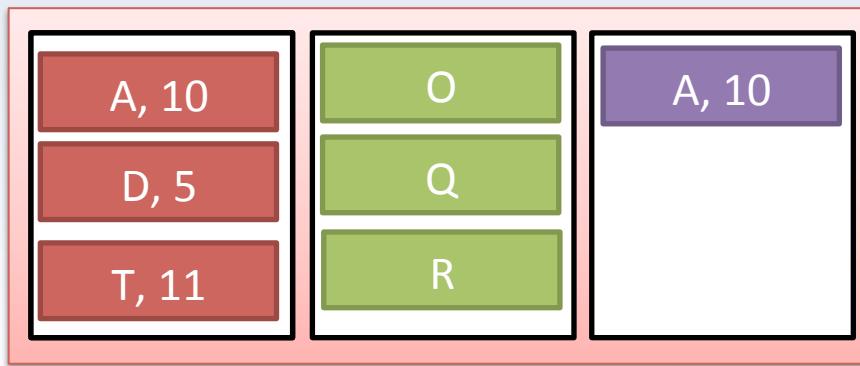
# One Pass Sort-Join



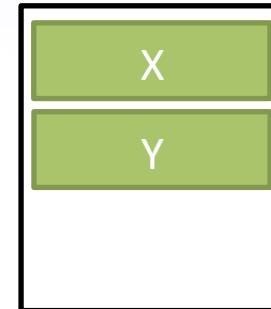
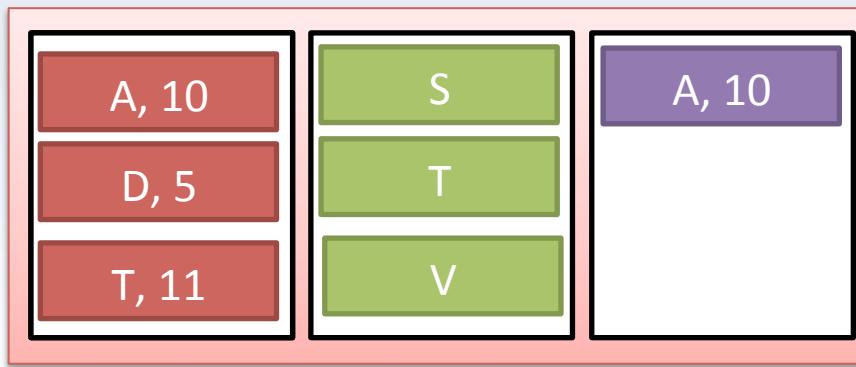
# One Pass Sort-Join



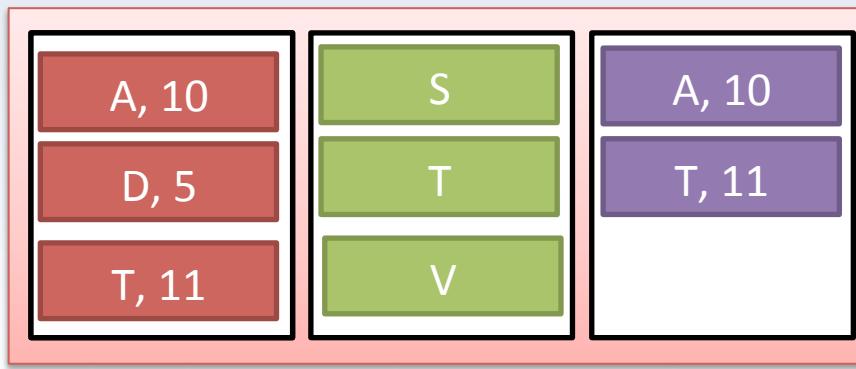
# One Pass Sort-Join



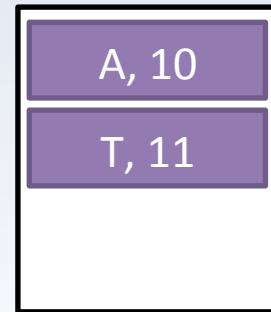
# One Pass Sort-Join

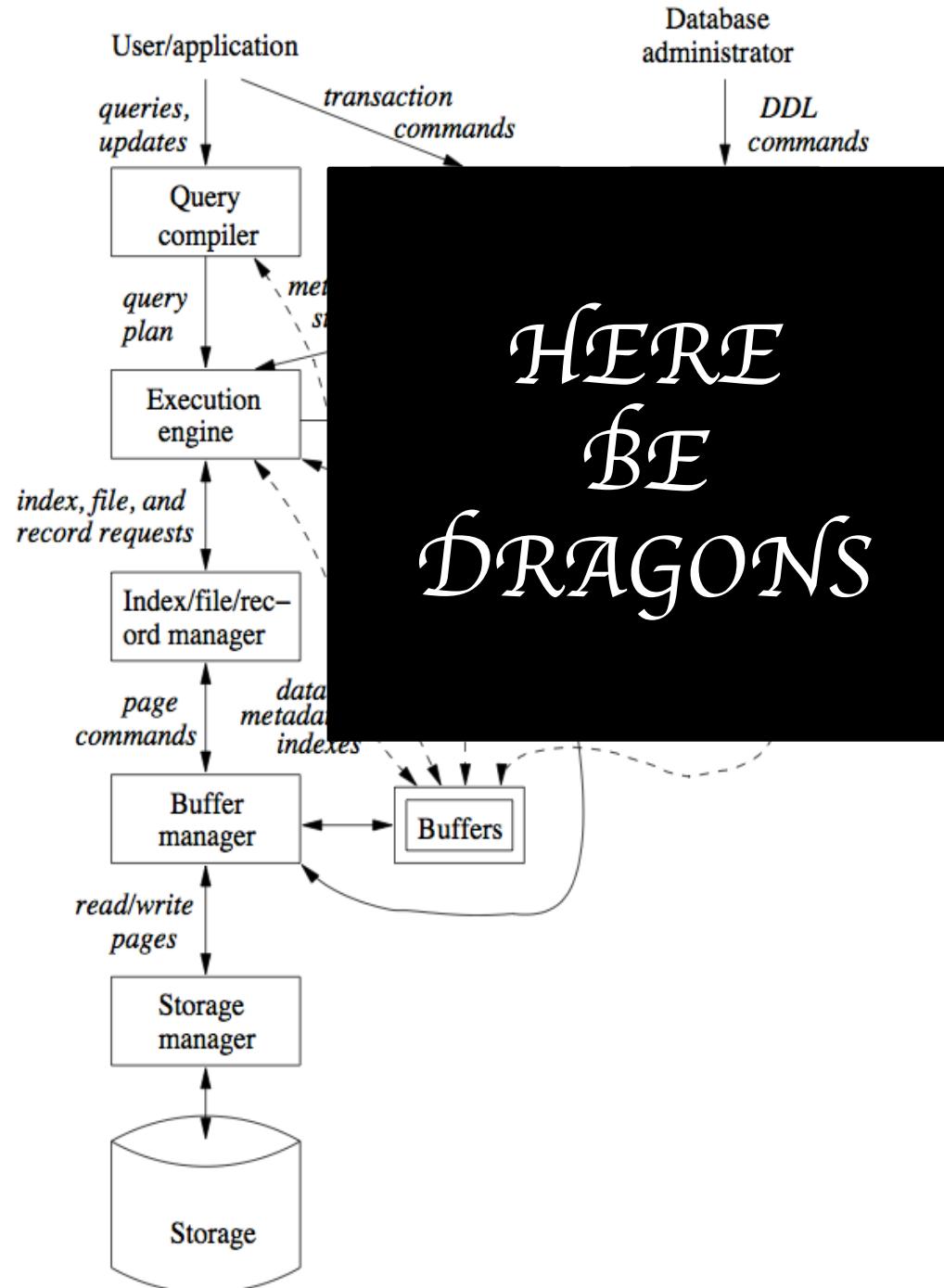


# One Pass Sort-Join



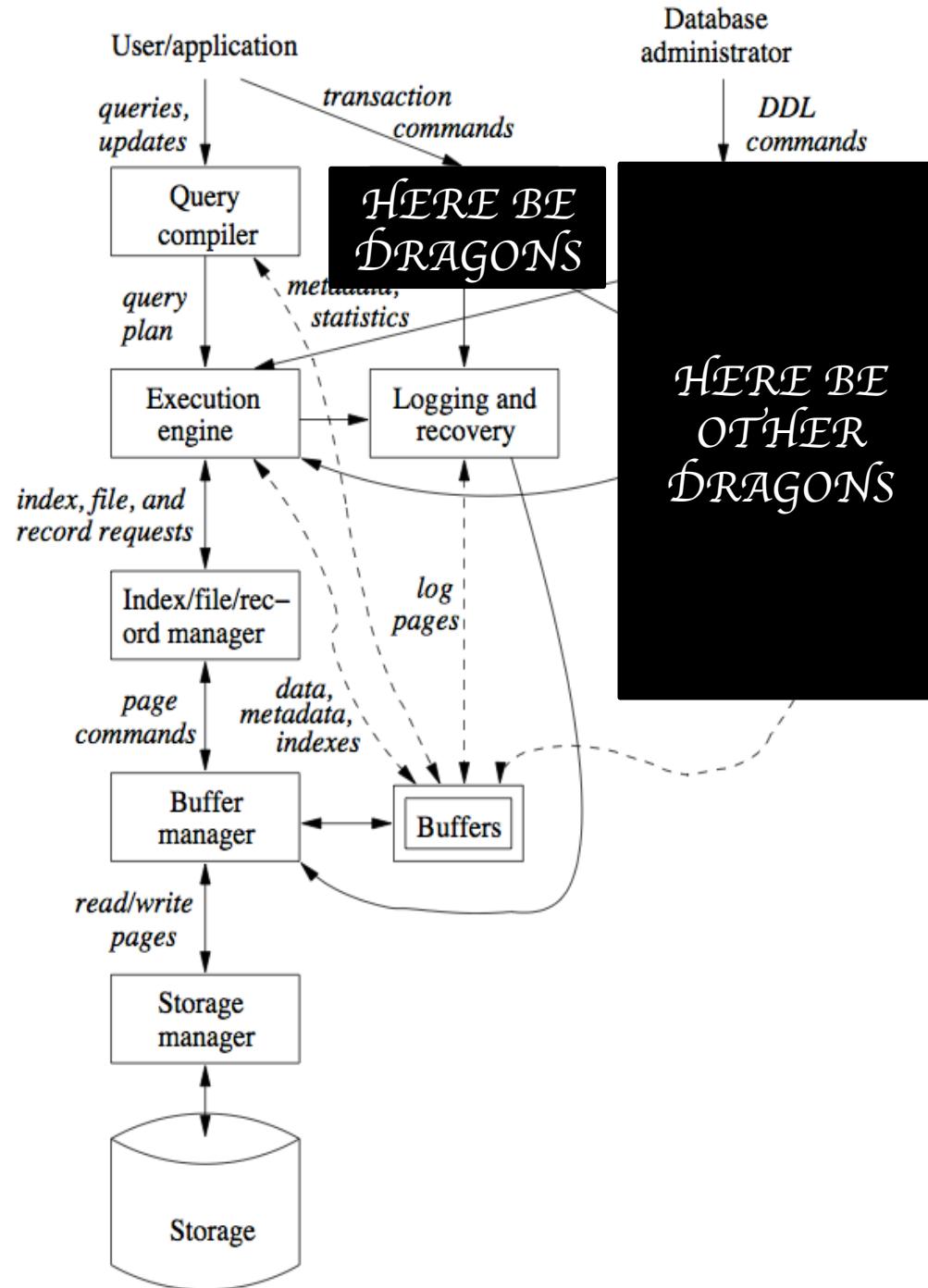
# Output







illinois.edu



# Logging and Recovery

- System failure or error: return the database to a valid state
  - transactions that are complete should be ***durable*** (written permanently to disk)
  - transactions that are incomplete should be ***aborted*** (appear to never have happened)



# Logging and Recovery

- ACID requirements:
  - atomicity
  - consistency
  - isolation
  - durability



# Logging and Recovery

- ACID requirements:
  - **atomicity** (Interruptions, ROLLBACK)
  - **consistency** (Constraints/Triggers)
  - isolation
  - **durability**



# Transaction Primitives

- $\text{INPUT}(X)$  - copy  $X$  into memory buffer
- $\text{READ}(X,t)$  - assign  $X$  to variable  $t$
- $\text{WRITE}(X,t)$  - copy value of  $t$  into  $X$
- $\text{OUTPUT}(X)$  - write memory buffer containing  $X$  to disk



# Example

Constraint:  $A=B$

Transaction:

$$A=A^*2$$
$$B=B^*2$$


# Example

Constraint:  $A=B$

Transaction:

$A=A^*2$

$B=B^*2$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t=t^* 2$	16	8		8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t = t^* 2$	16	8		8	8
WRITE(A,t)	16	16		8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t = t^* 2$	16	8		8	8
WRITE(A,t)	16	16		8	8
READ(B,t)	8	16	8	8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t=t^* 2$	16	8		8	8
WRITE(A,t)	16	16		8	8
READ(B,t)	8	16	8	8	8
$t=t^* 2$	16	16	8	8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t=t^*2$	16	8		8	8
WRITE(A,t)	16	16		8	8
READ(B,t)	8	16	8	8	8
$t=t^*2$	16	16	8	8	8
WRITE(B,t)	16	16	16	8	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t=t^*2$	16	8		8	8
WRITE(A,t)	16	16		8	8
READ(B,t)	8	16	8	8	8
$t=t^*2$	16	16	8	8	8
WRITE(B,t)	16	16	16	8	8
OUTPUT(A)	16	16	16	16	8



# Example

Constraint:  $A=B$

Transaction:

$$A = A^* 2$$

$$B = B^* 2$$

Action	t	Mem A	Mem B	Disk A	Disk B
READ(A,t)	8	8		8	8
$t=t^*2$	16	8		8	8
WRITE(A,t)	16	16		8	8
READ(B,t)	8	16	8	8	8
$t=t^*2$	16	16	8	8	8
WRITE(B,t)	16	16	16	8	8
OUTPUT(A)	16	16	16	16	8
OUTPUT(B)	16	16	16	16	16



# Undo Logging

- Main idea: keep track of all values we change
- When recovering, restore the old values for incomplete transactions



# Logging

- Sequentially write entries to disk:
  - $\langle \text{START } T \rangle$  - transaction started
  - $\langle \text{COMMIT } T \rangle$  - transaction successful
  - $\langle \text{ABORT } T \rangle$  - transaction unsuccessful
  - $\langle T, X, v \rangle$  - transaction changed element X from value v to some new value



# Undo Logging

- Rules:
  1. Add  $\langle T, X, v \rangle$  to log ***before*** new value of X is written to disk
  2. Add  $\langle \text{COMMIT } T \rangle$  ***after all*** results for T have been written to disk



# Example

Action	t	Mem A	Mem B	Disk A	Disk B	Log
						<START T>
READ(A,t)	8	8		8	8	
t=t*2	16	8		8	8	
WRITE(A,t)	16	16		8	8	<T,A,8>
READ(B,t)	8	16	8	8	8	
t=t*2	16	16	8	8	8	
WRITE(B,t)	16	16	16	8	8	<T,B,8>
FLUSH LOG						
OUTPUT(A)	16	16	16	16	8	
OUTPUT(B)	16	16	16	16	16	
						<COMMIT T>
FLUSH LOG						



# Recovery

- If we see  $\langle \text{START } T \rangle$  without corresponding  $\langle \text{COMMIT } T \rangle$ , must undo the transaction
- When the system fails, scan log backward
  - If we see  $\langle \text{COMMIT } F \rangle$ , ignore transaction F
  - For every other transaction T, when we see  $\langle T, X, v \rangle$ , write v to X



# Example

$\langle \text{START } T_1 \rangle$

$\langle T_1, A, 5 \rangle$

$\langle T_1, B, 3 \rangle$

$\langle \text{START } T_2 \rangle$

$\langle T_1, C, 7 \rangle$

$\langle T_2, X, 5 \rangle$

$\langle \text{COMMIT } T_2 \rangle$

$\langle \text{START } T_3 \rangle$

$\langle T_1, D, 1 \rangle$

$\langle T_3, A, 3 \rangle$

$\langle T_3, B, 2 \rangle$

$\langle \text{START } T_4 \rangle$

$\langle \text{COMMIT } T_3 \rangle$

$\langle T_4, X, 3 \rangle$



# Next time...

- We'll look at
  - Ways to avoid scanning entire log
  - Redo logging
  - Undo/Redo logging
  - Introduce transaction management

