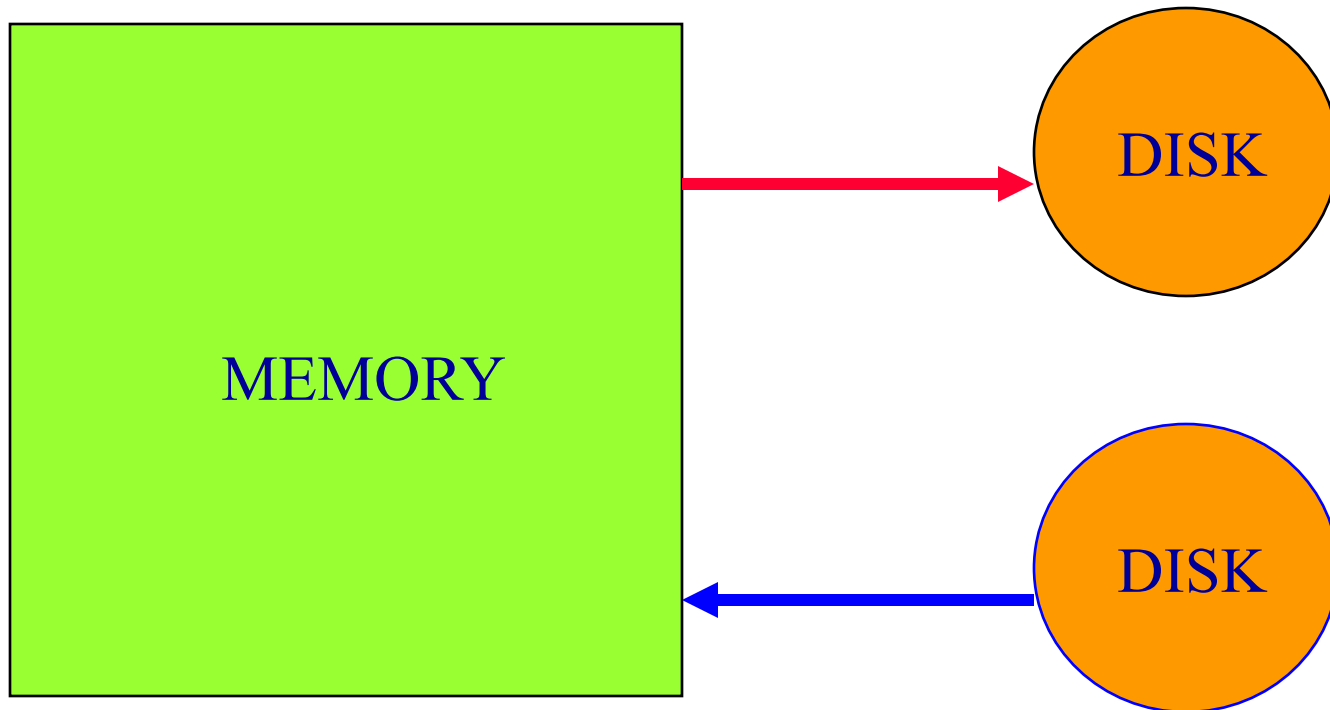
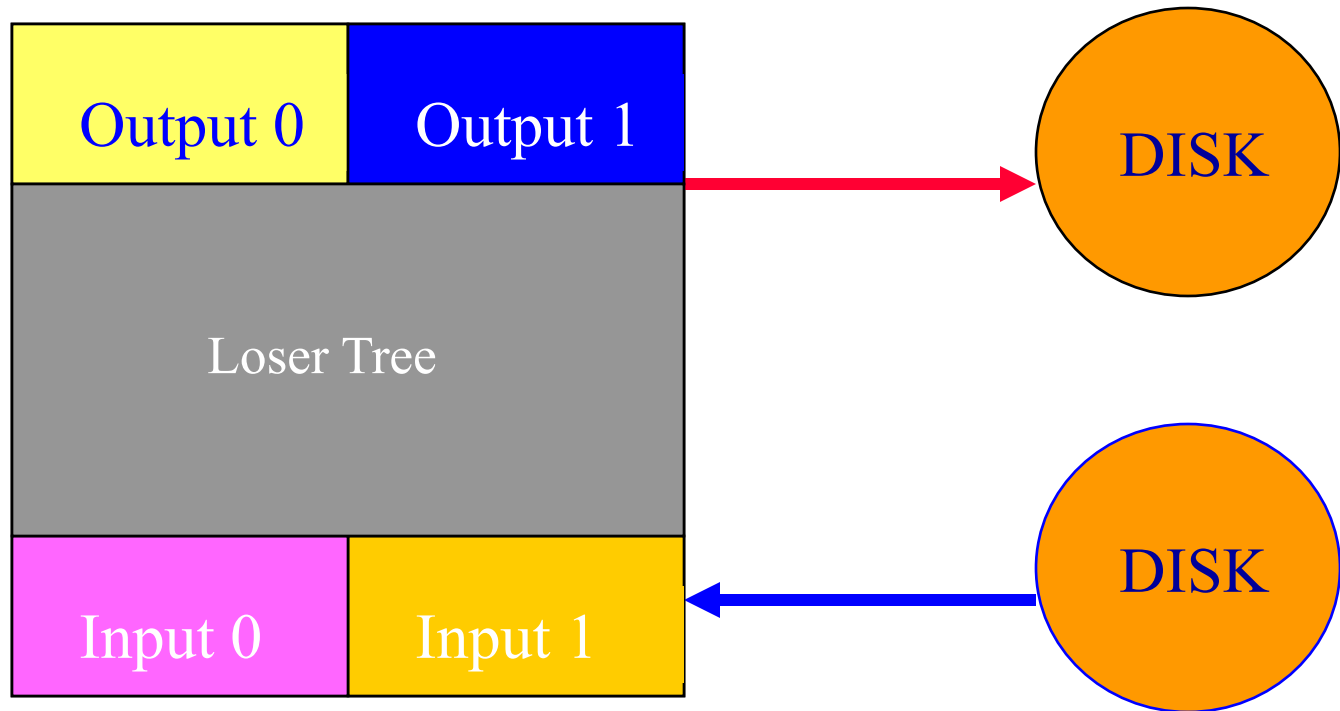


# Improve Run Generation

- Overlap input, output, and internal CPU work.
- Reduce the number of runs (equivalently, increase average run length).

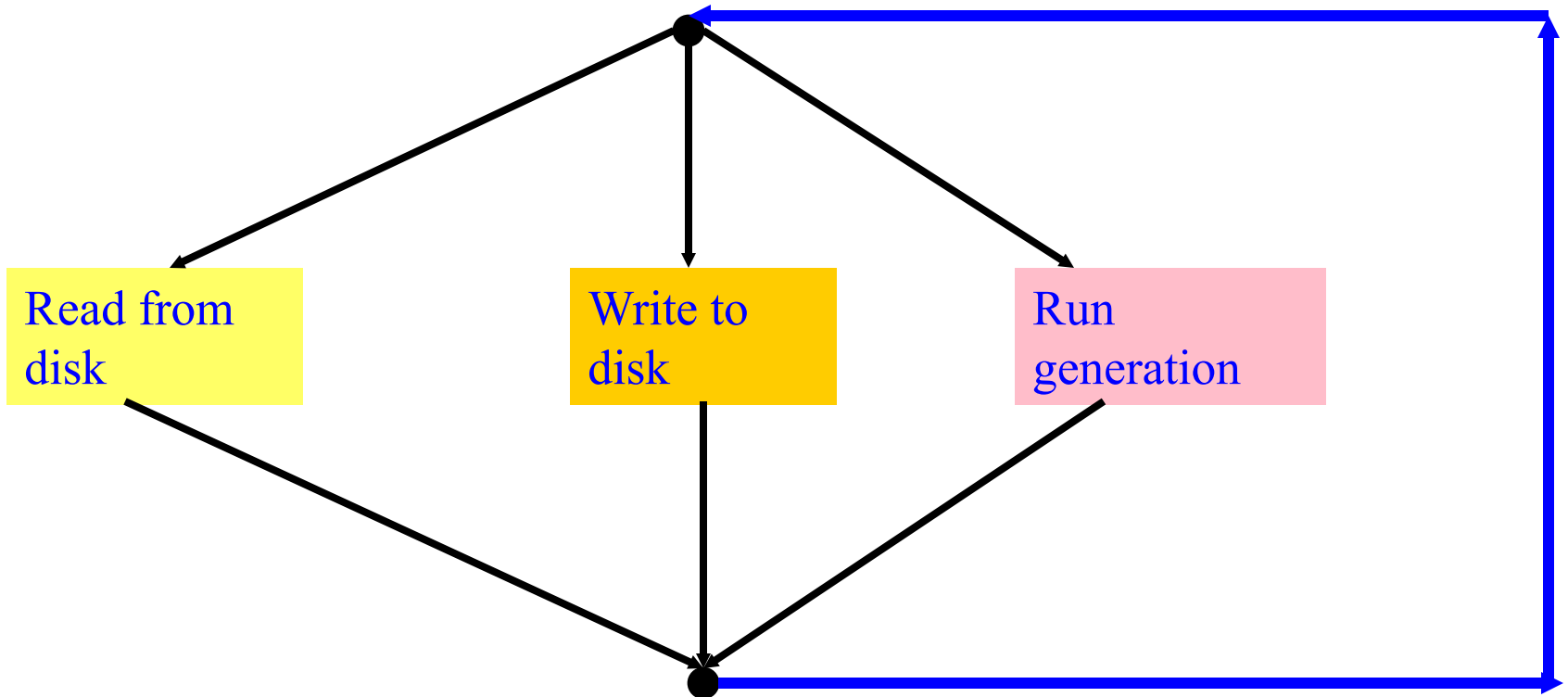


# New Strategy



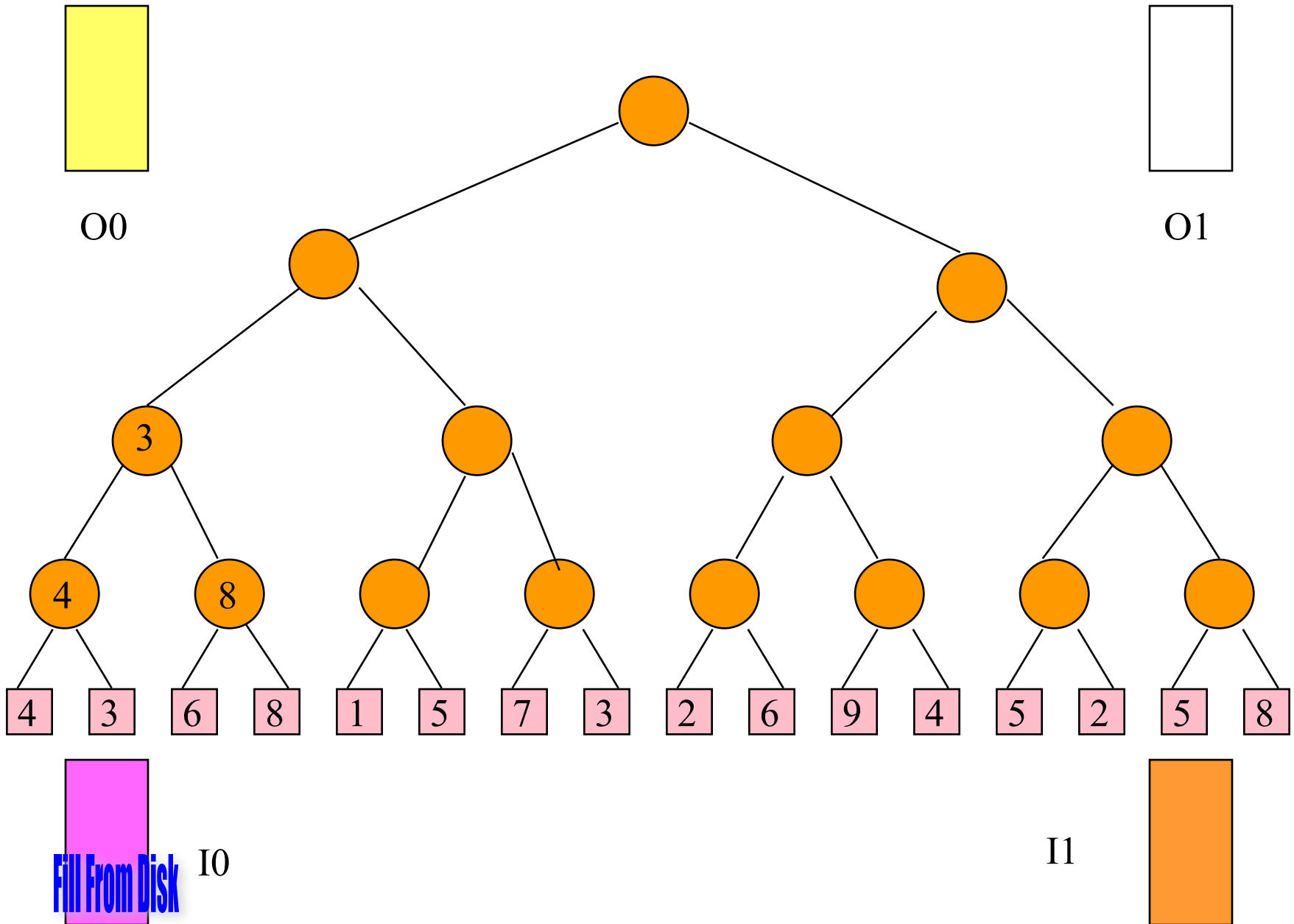
- Use 2 input and 2 output buffers.
- Rest of memory is used for a min loser tree.
- Actually, 3 buffers adequate.

# Steady State Operation

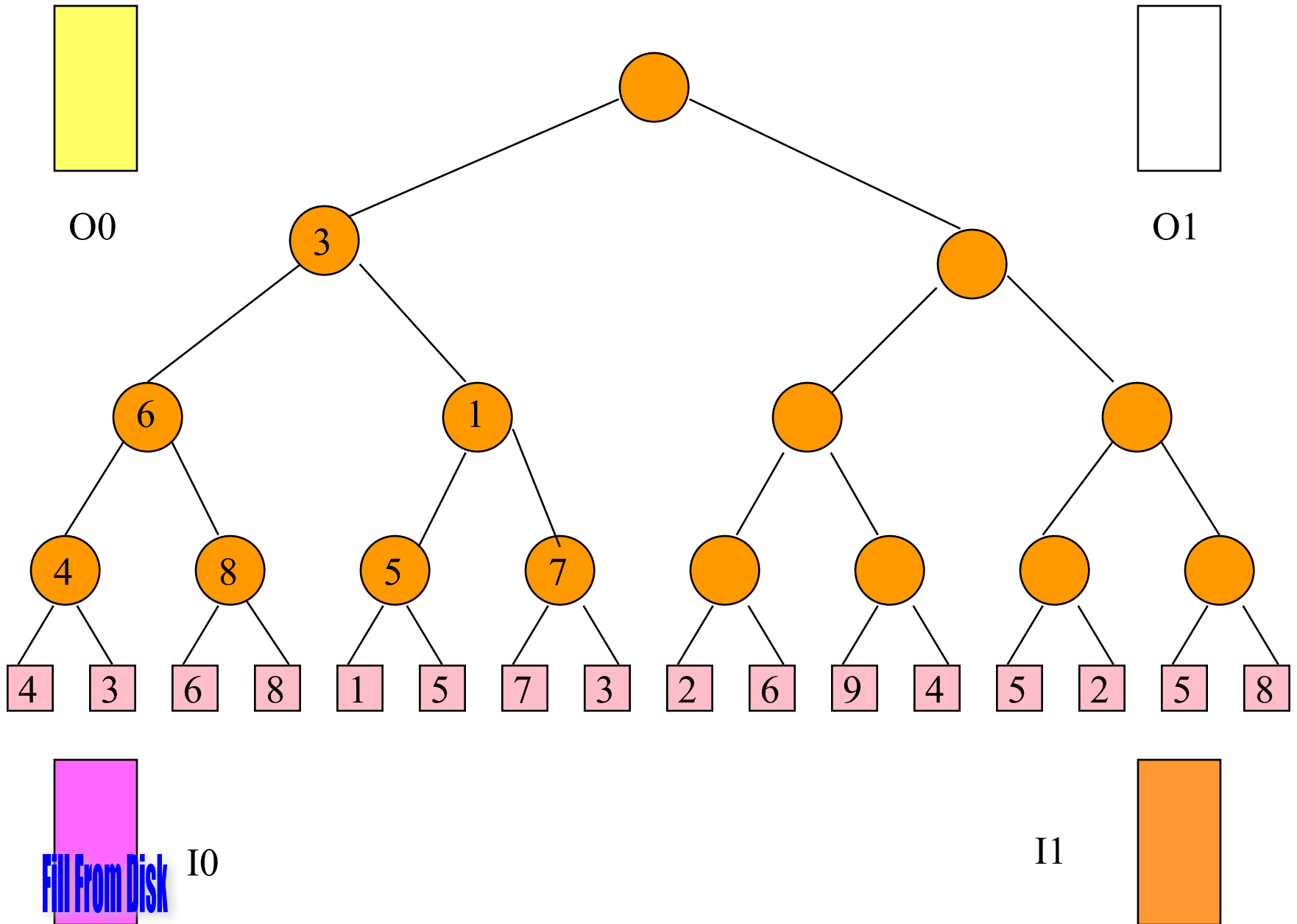


- Synchronization is done when the active input buffer gets empty (the active output buffer will be full at this time).

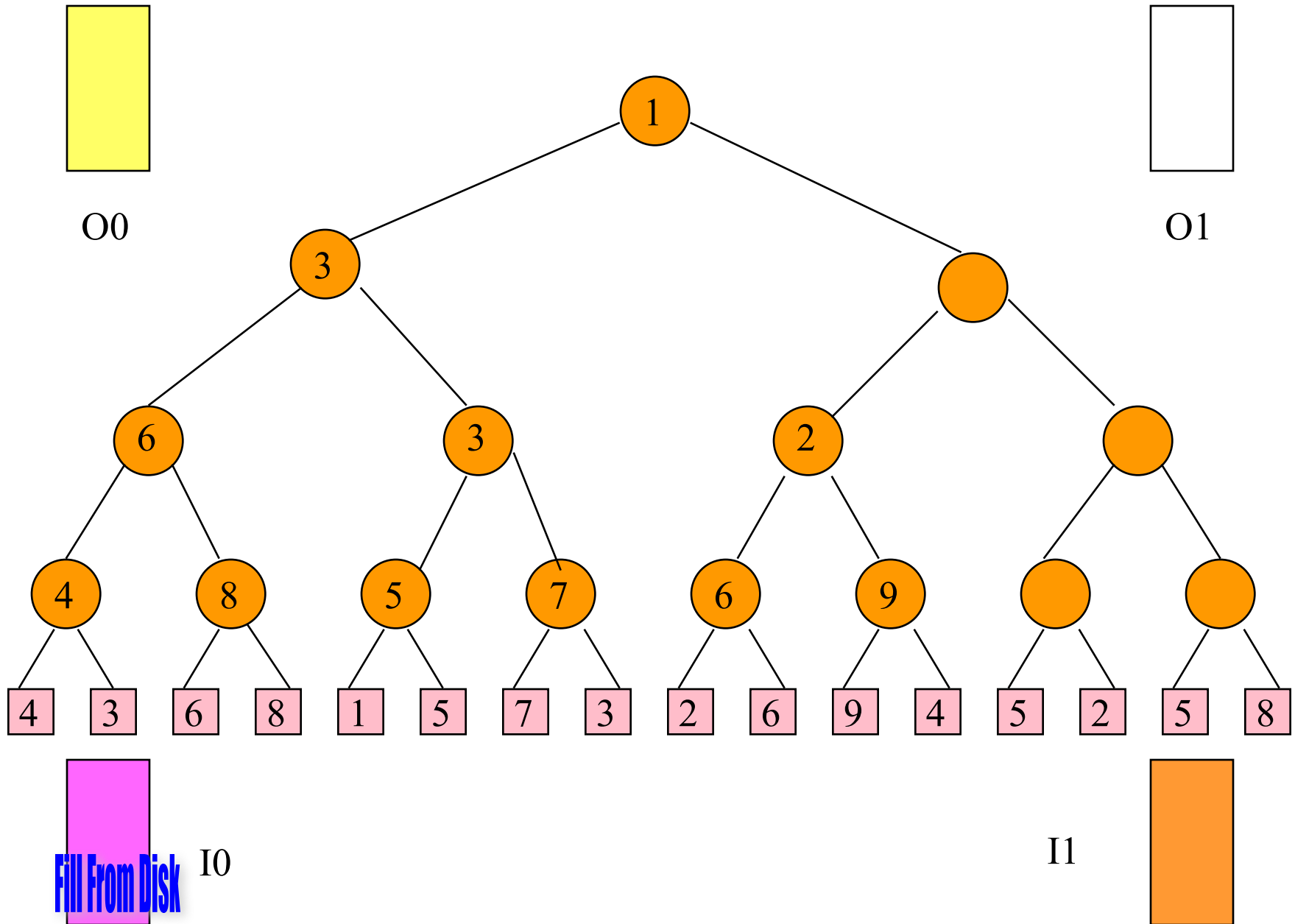
# Initialize



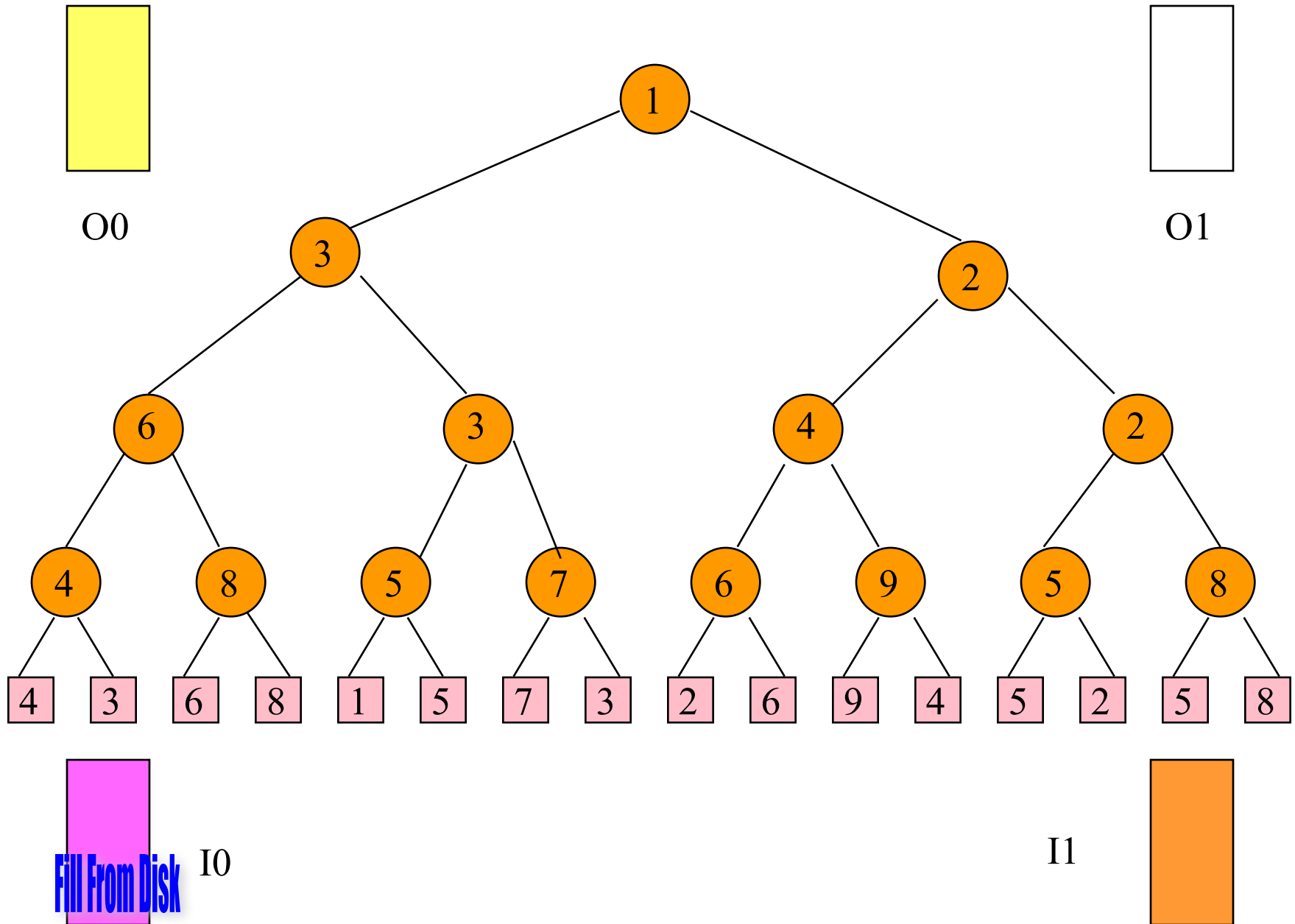
# Initialize



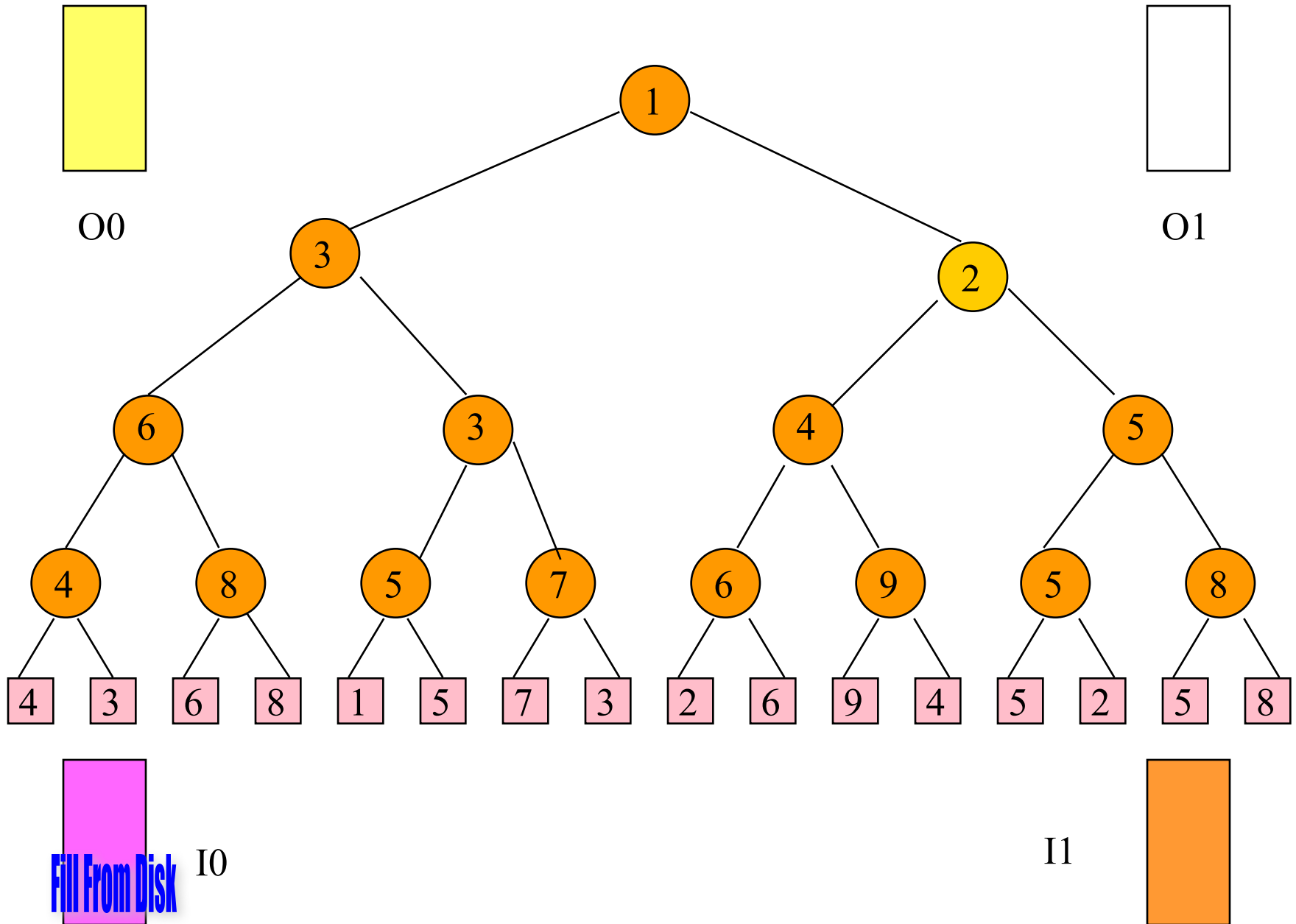
# Initialize



# Initialize

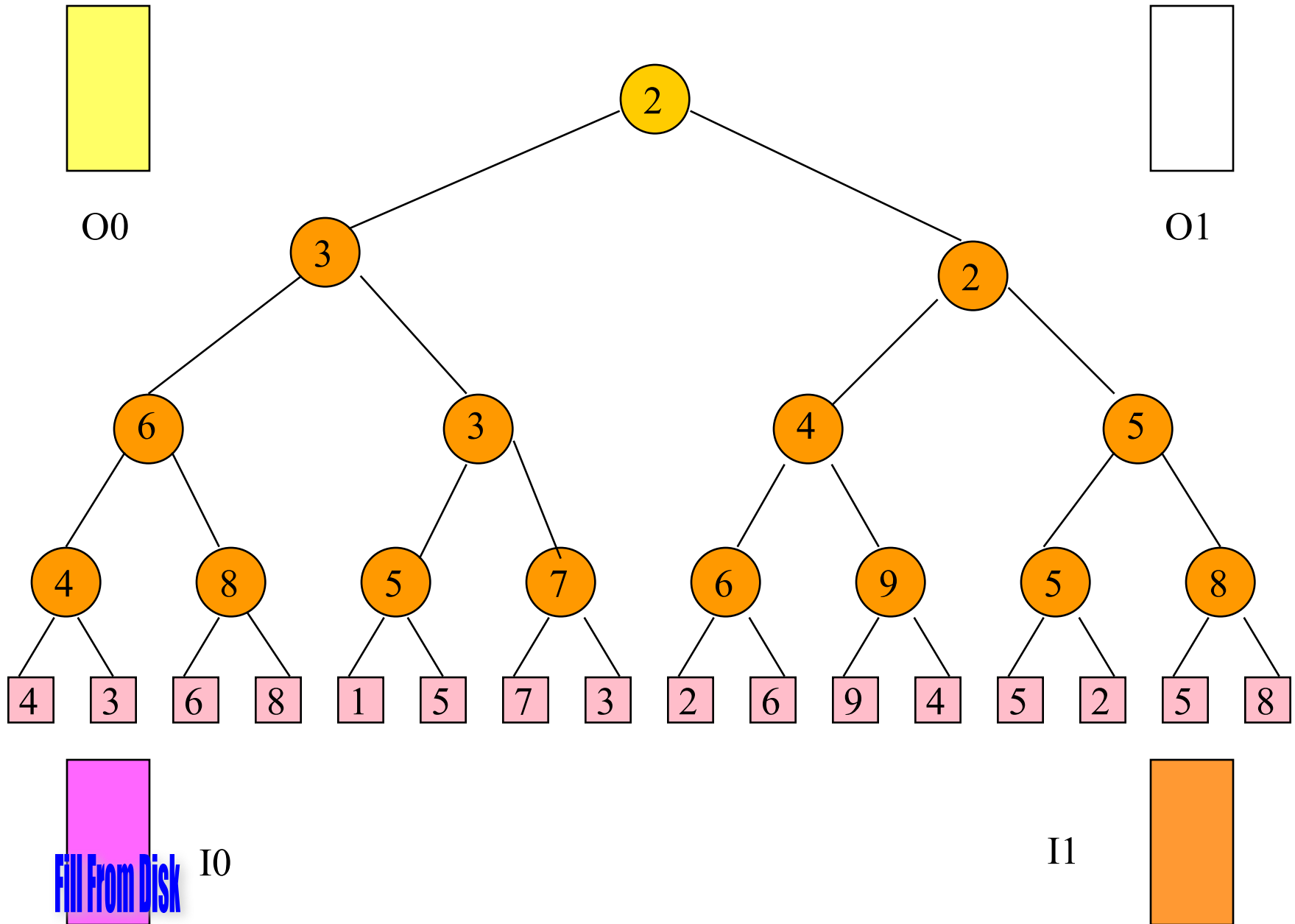


# Initialize



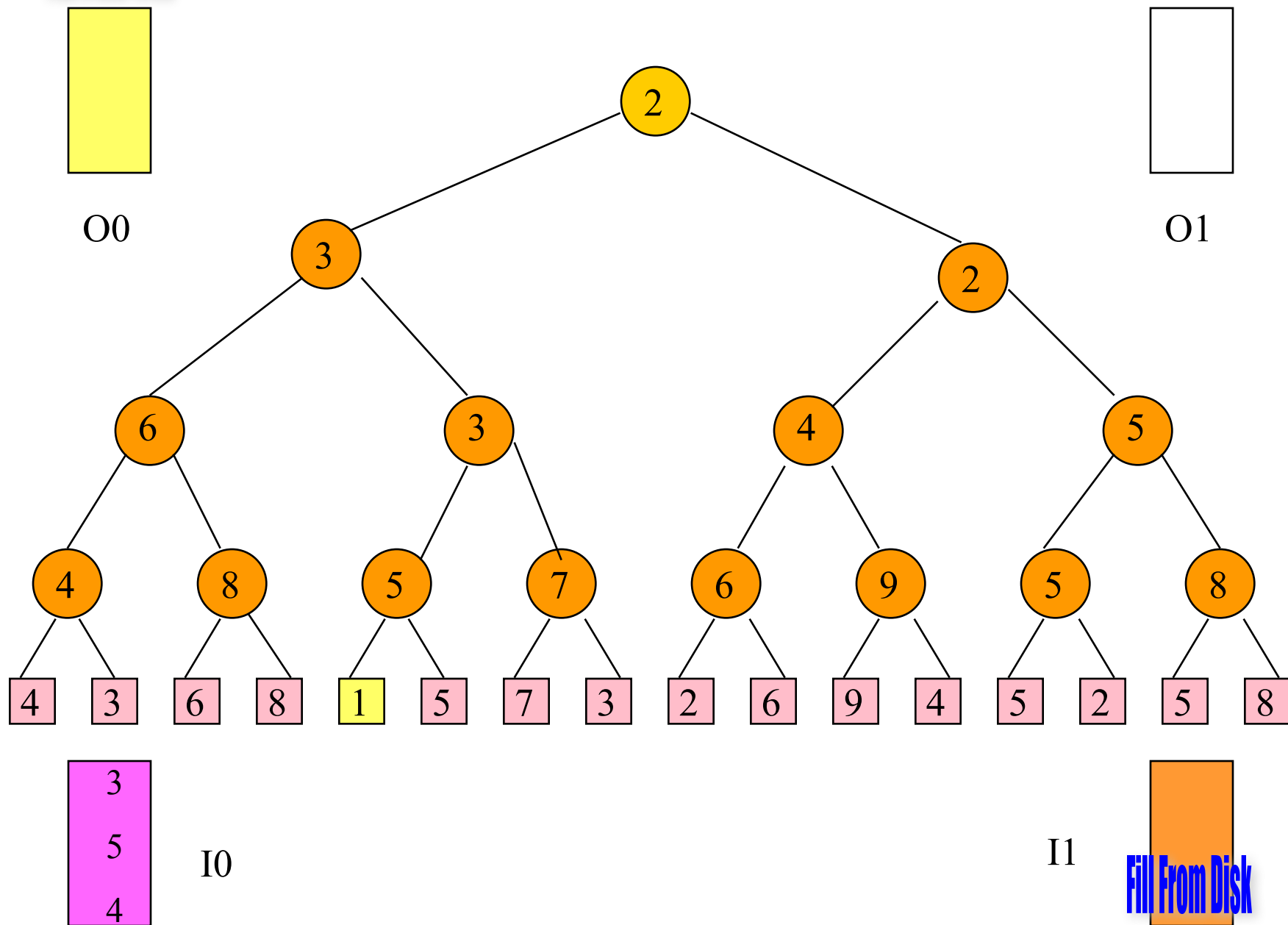


# Initialize

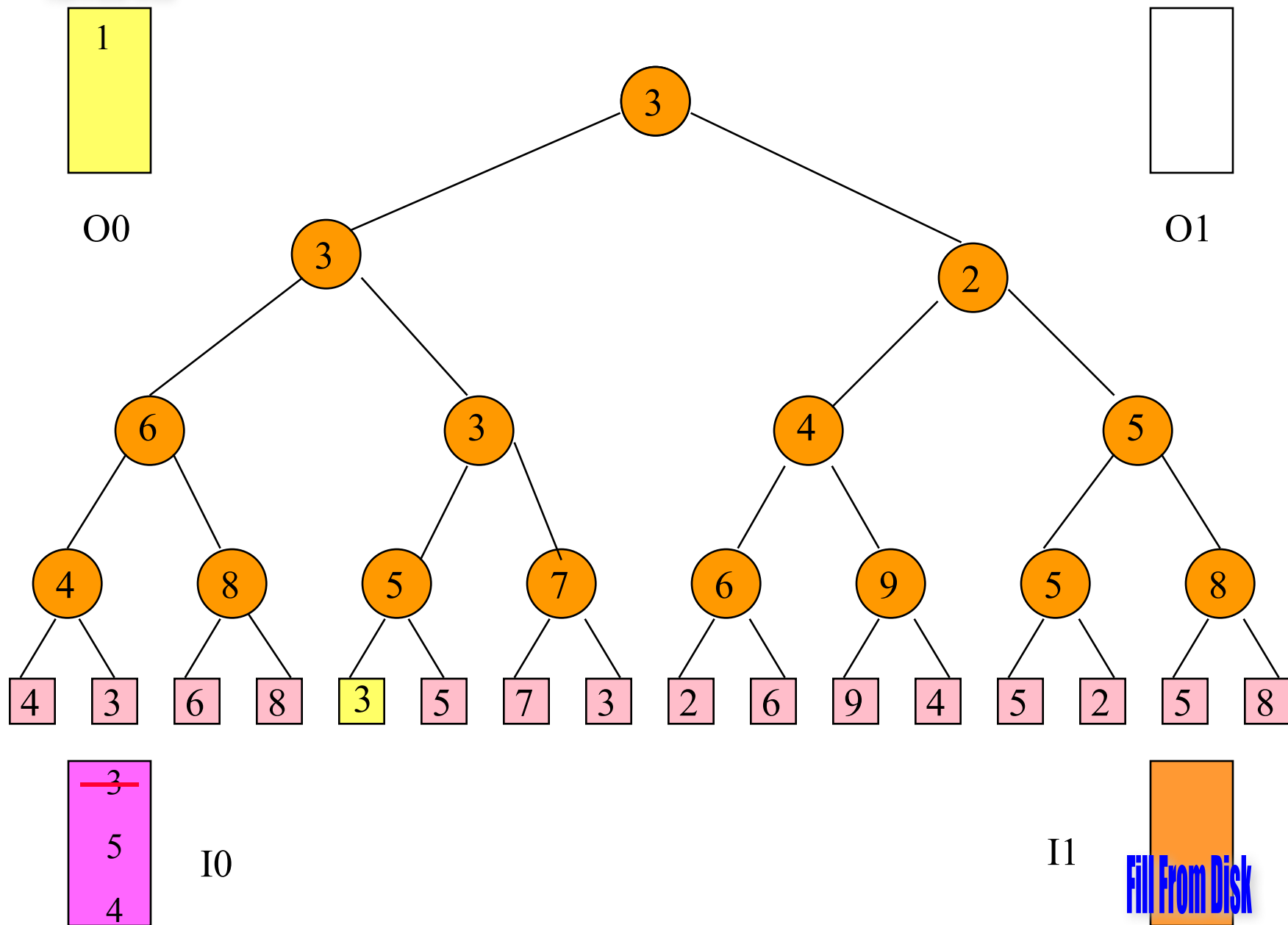


Fill From Tree

# Generate Run 1

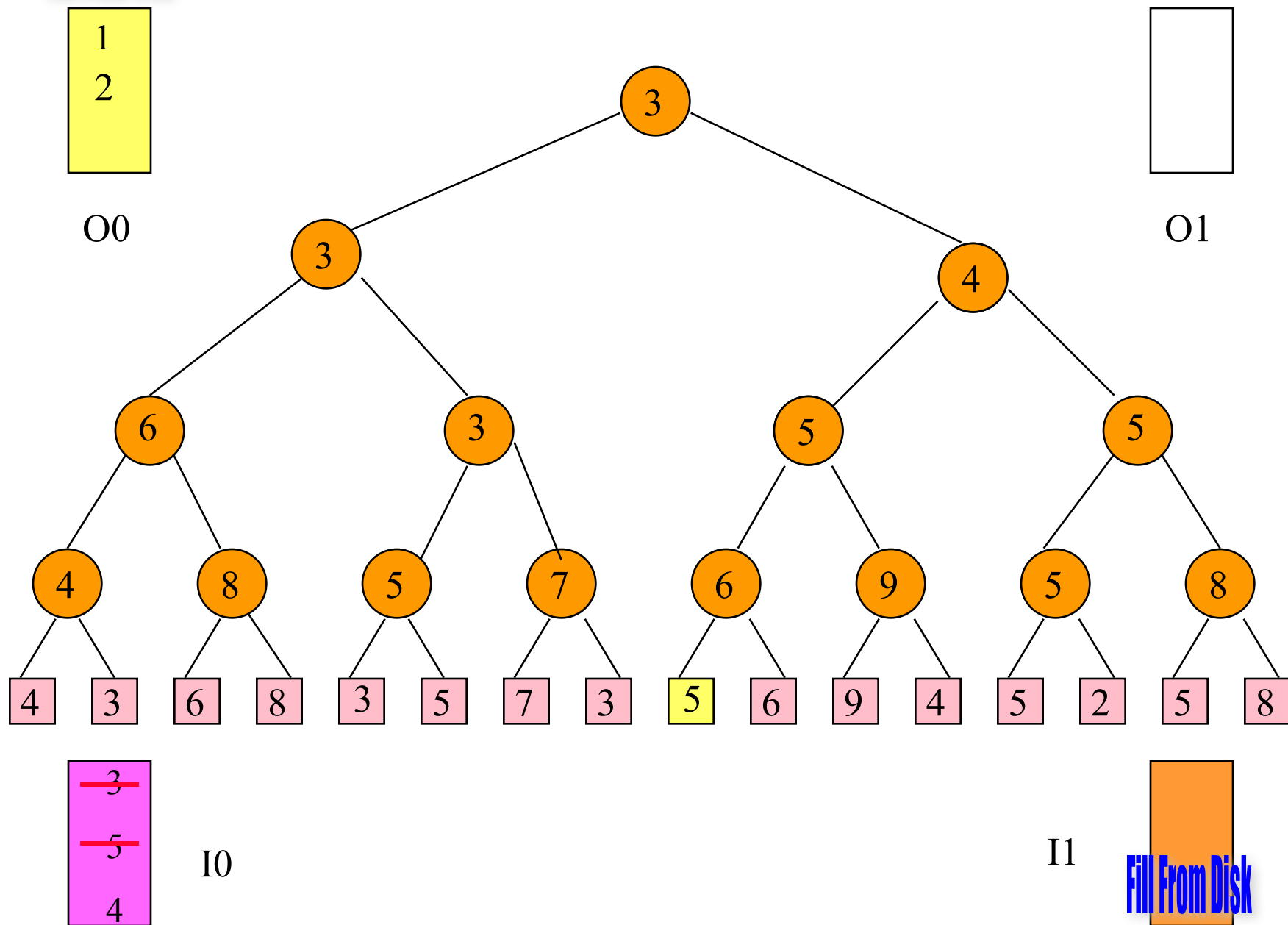


# Generate Run 1



Fill From Tree

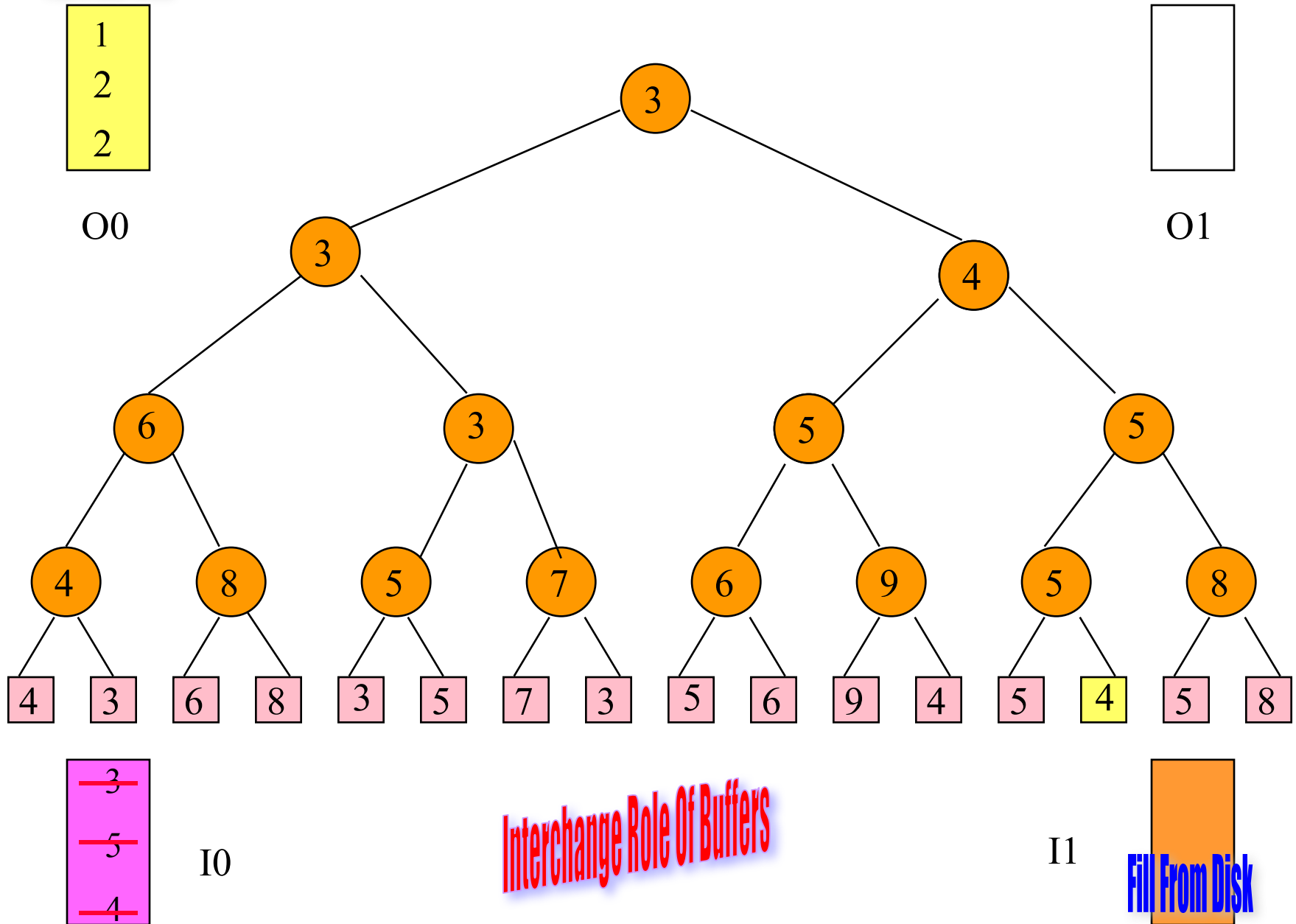
# Generate Run 1



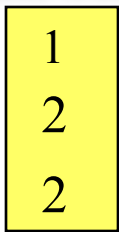
Fill From Disk

Fill From Tree

# Generate Run 1



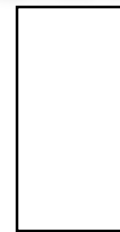
Write To Disk



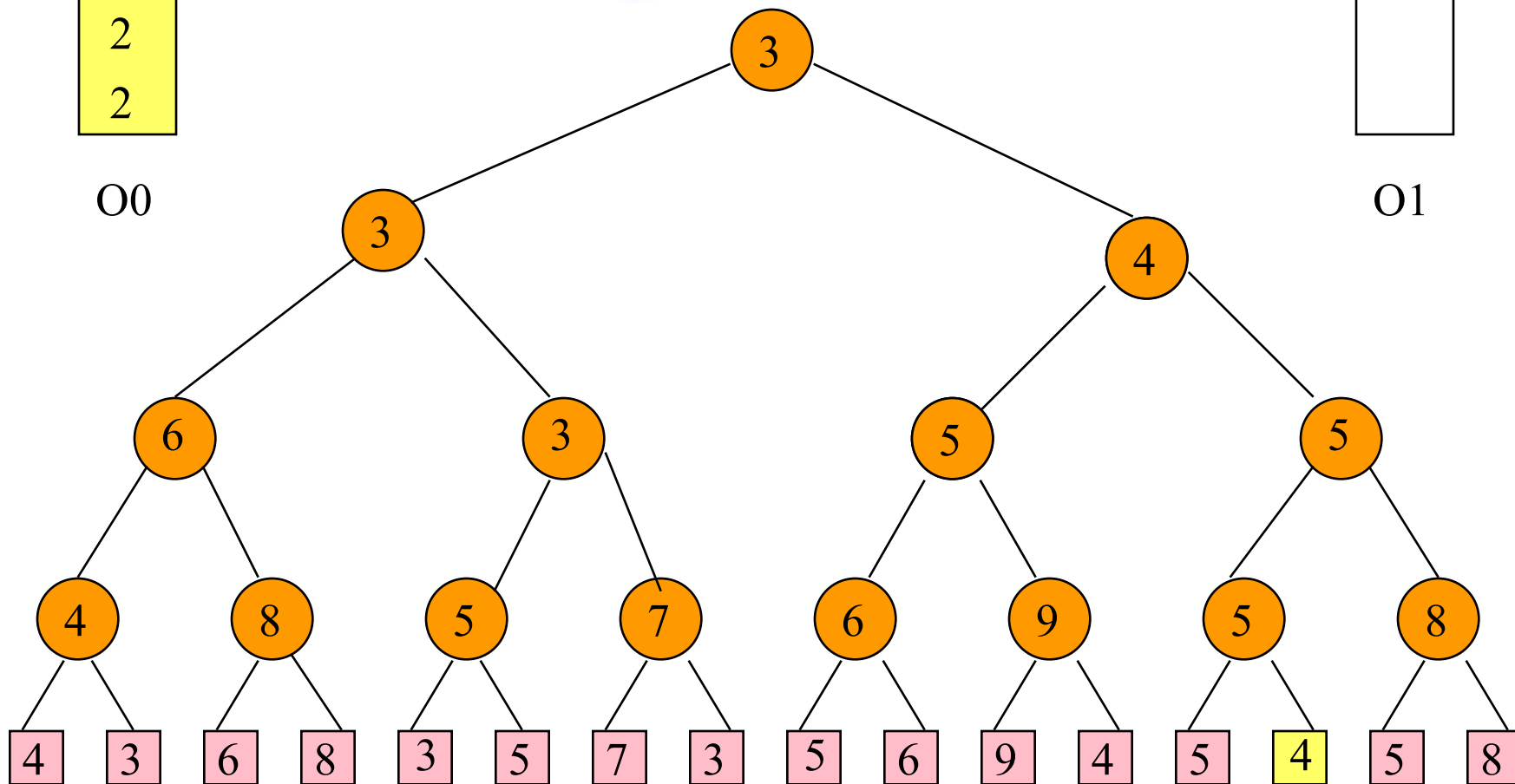
O0

Interchange Role Of Buffers

Fill From Tree



O1



Fill From Disk



I0

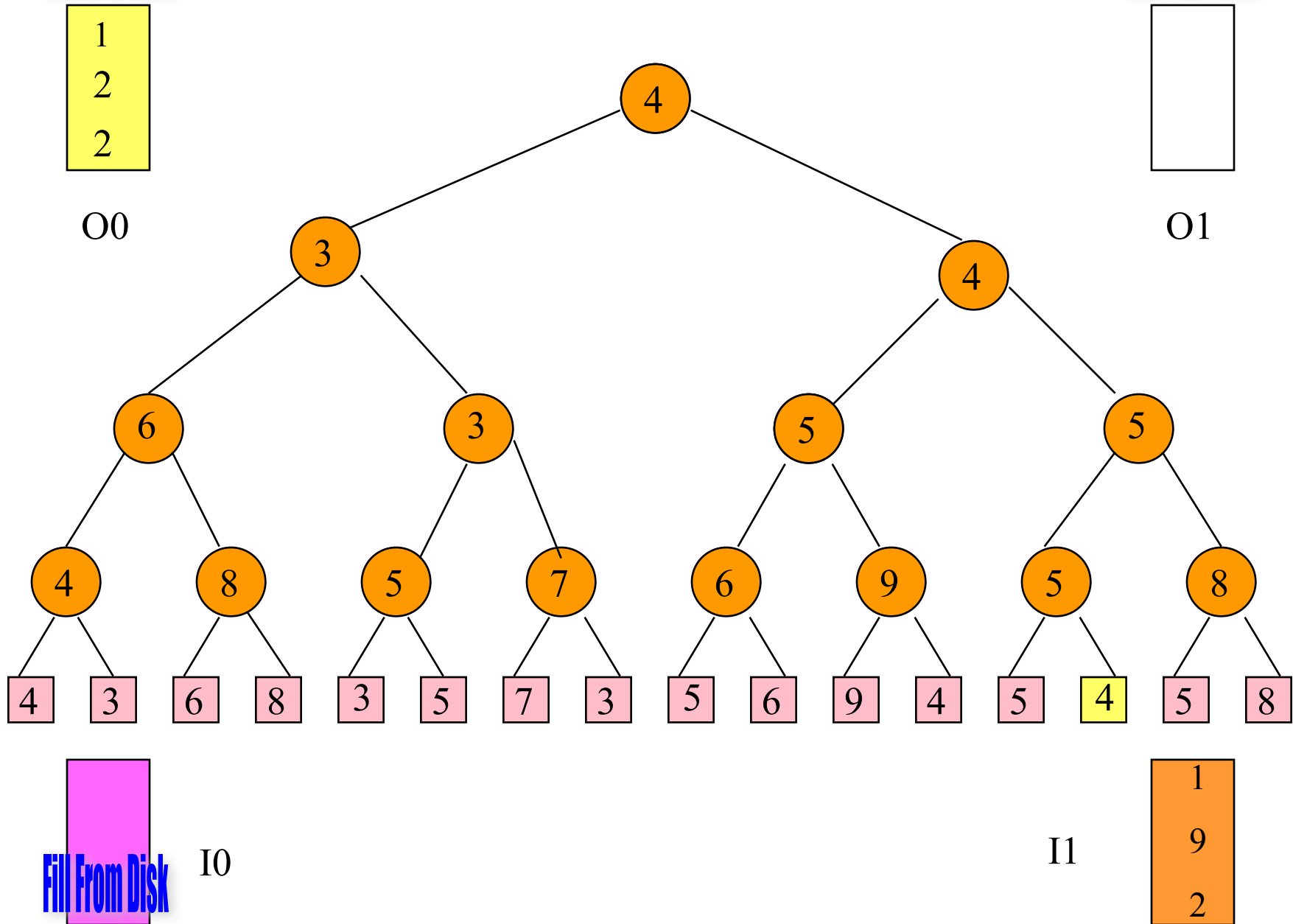
I1



## Write To Disk

# Continue With Run 1

## Fill From Tree



Write To Disk

1  
2  
2

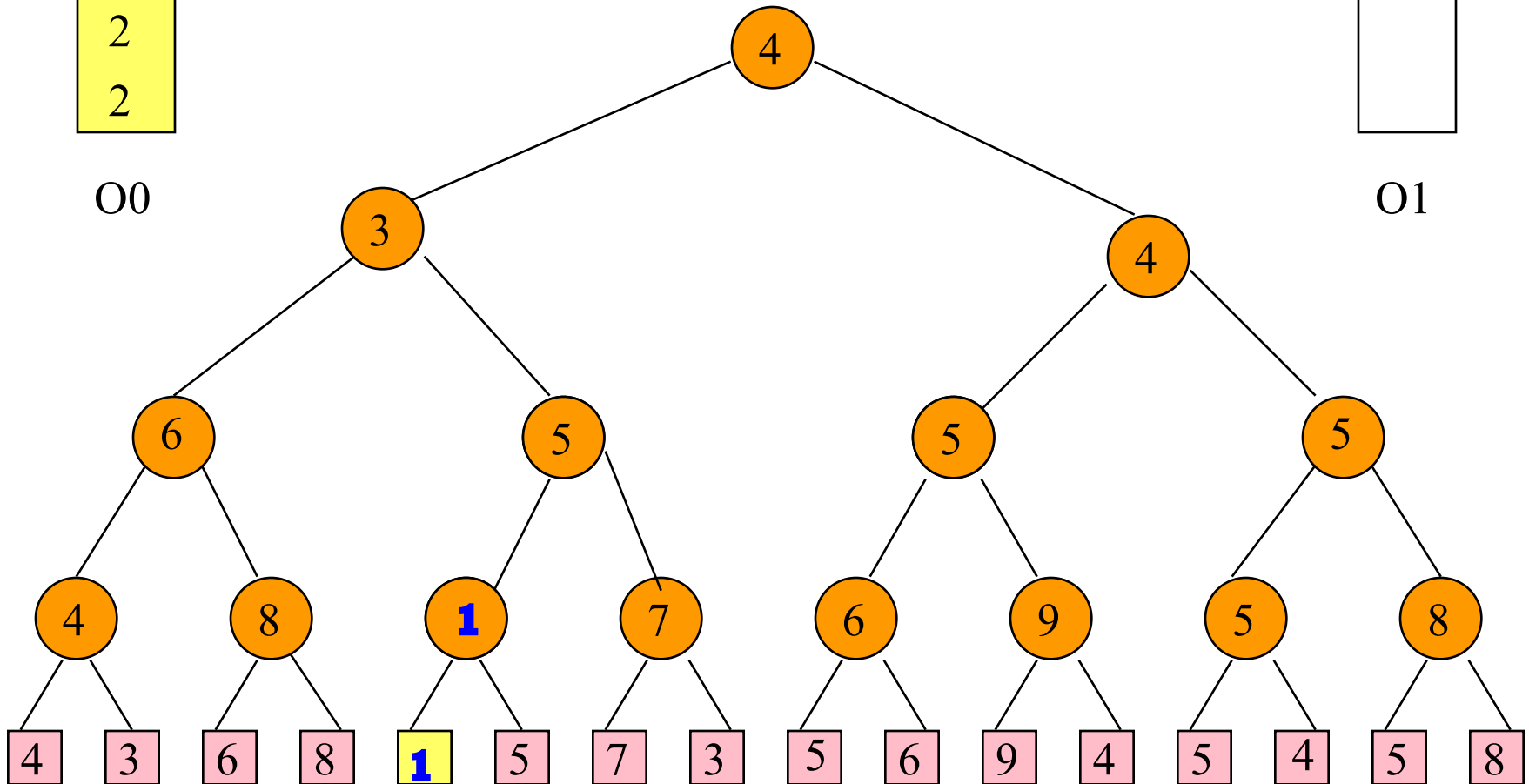
O0

# Continue With Run 1

Fill From Tree

3

O1



Fill From Disk

I0

I1

~~1~~  
9  
2



Write To Disk

1  
2  
2

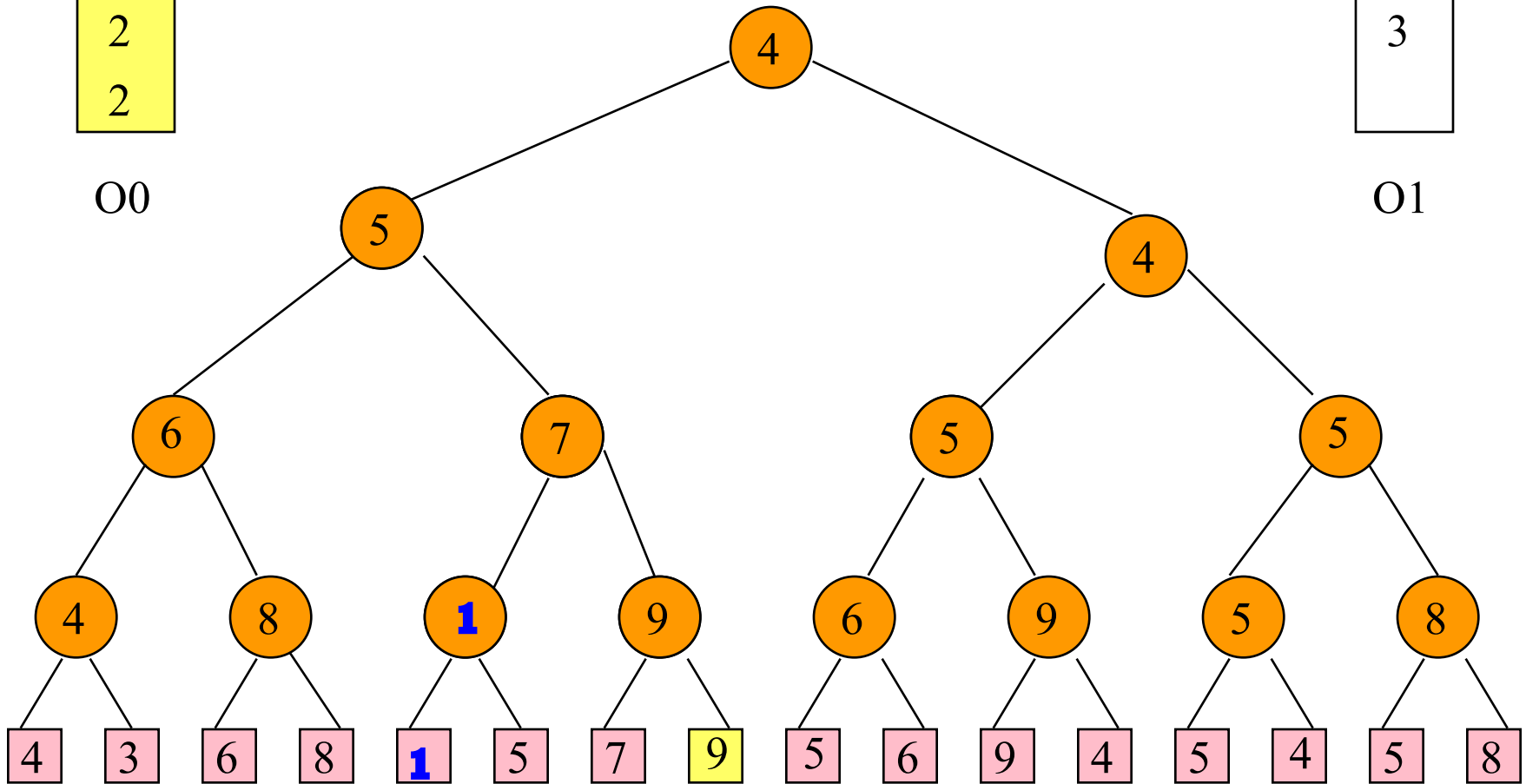
O0

# Continue With Run 1

Fill From Tree

3  
3

O1



Fill From Disk

I0

I1

~~1~~  
~~9~~  
2

Write To Disk

1  
2  
2

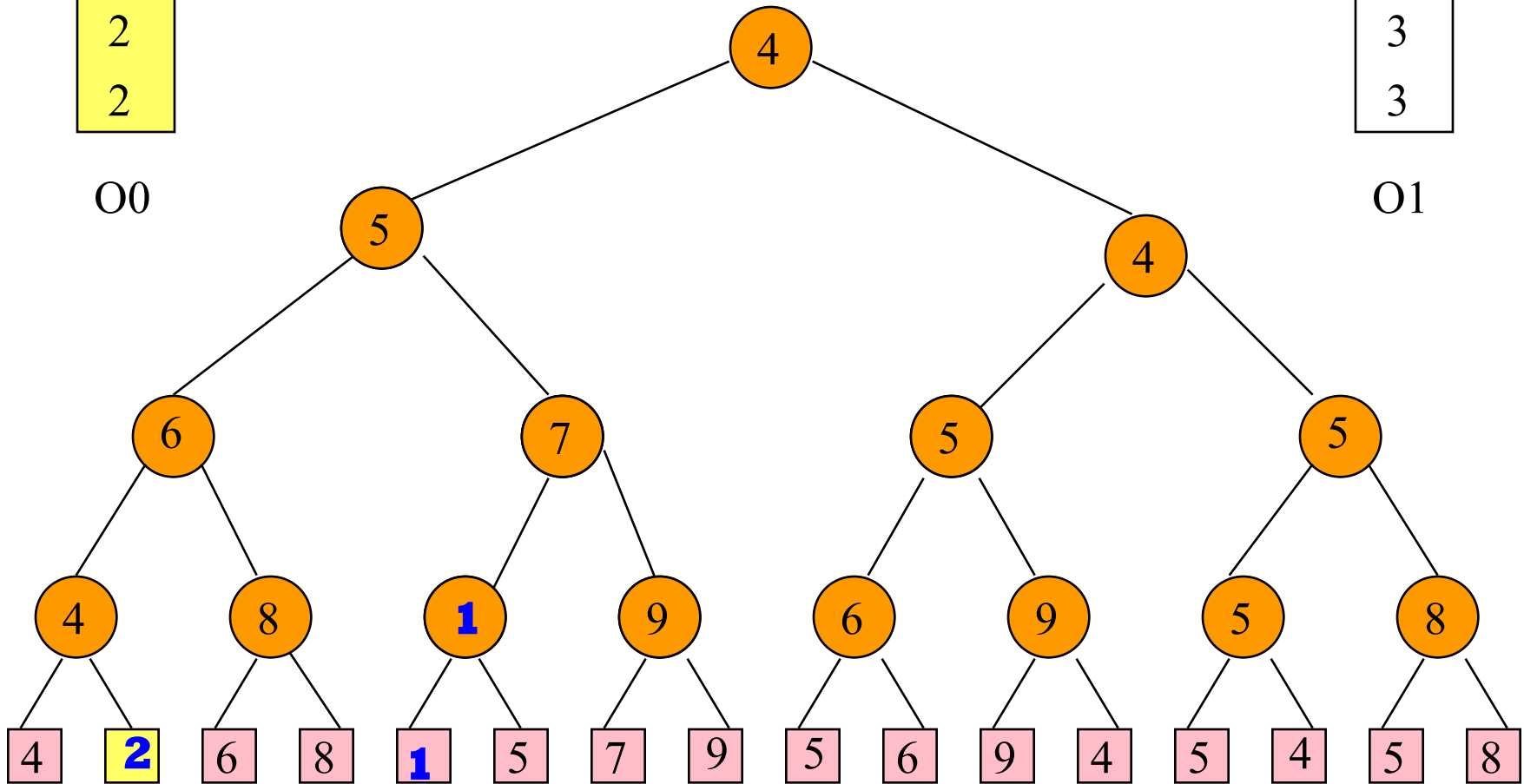
O0

# Continue With Run 1

Fill From Tree

3  
3  
3

O1



Fill From Disk

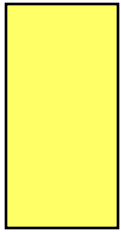
I0

Interchange Role Of Buffers

I1

~~1~~  
~~9~~  
~~2~~

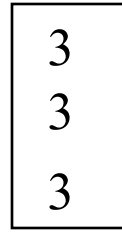
Fill From Tree



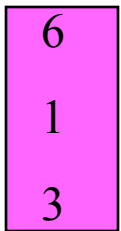
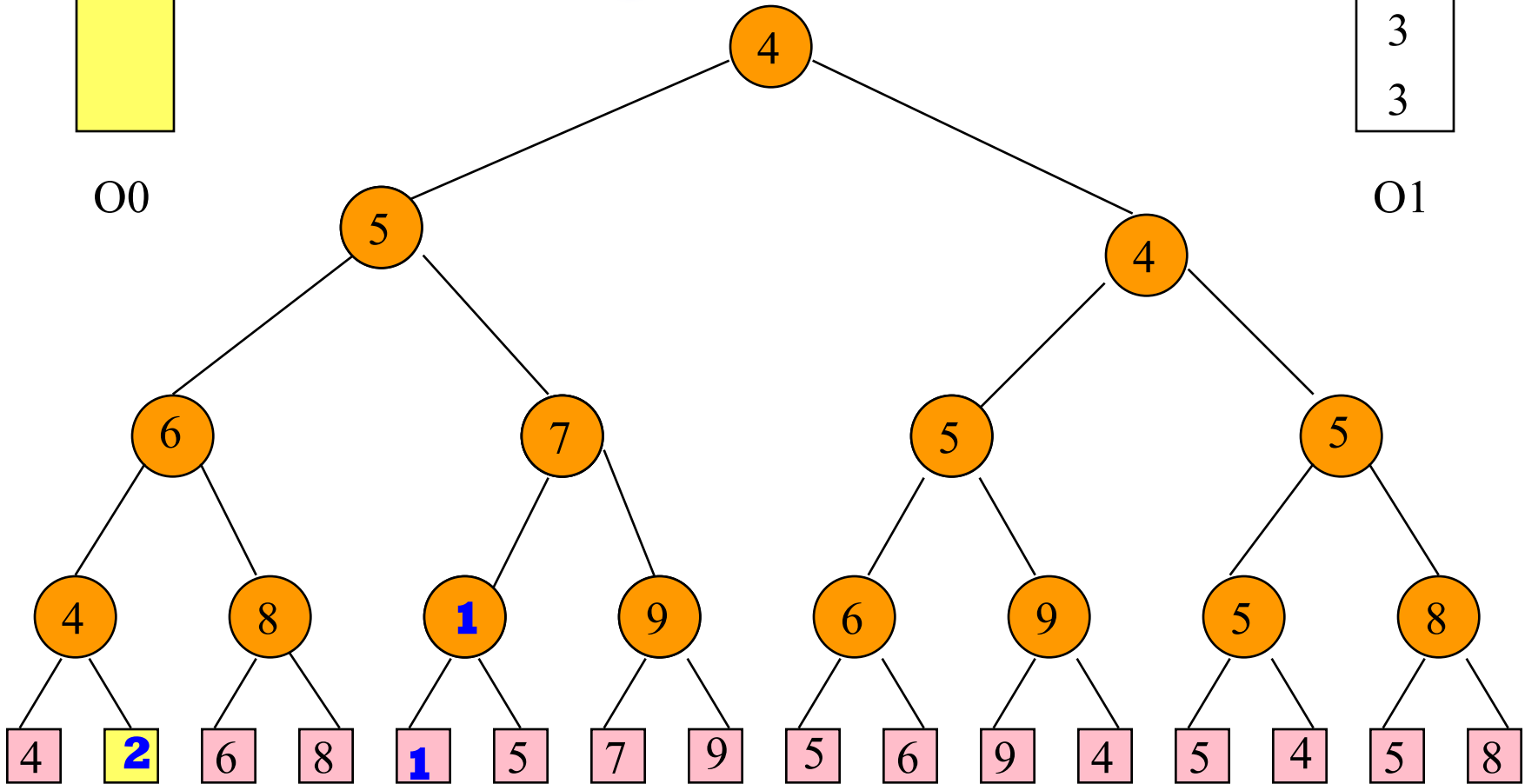
O0

Interchange Role Of Buffers

Write To Disk



O1



I0

I1

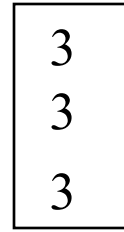
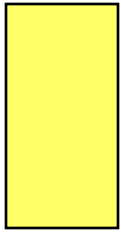


Fill From Disk

Fill From Tree

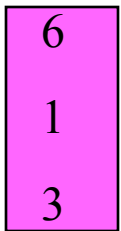
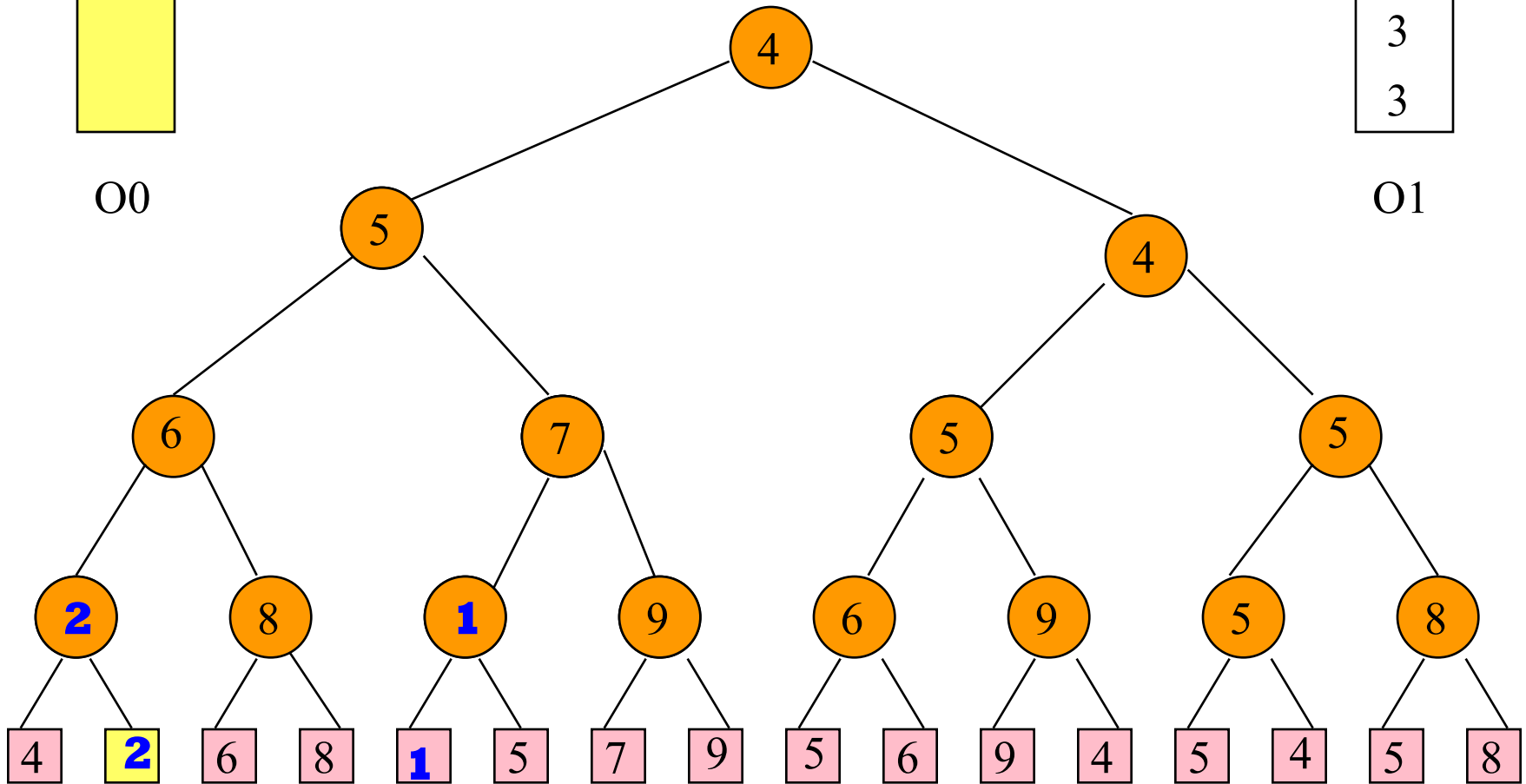
# Continue With Run 1

Write To Disk



O0

O1



I0

I1

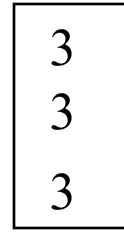
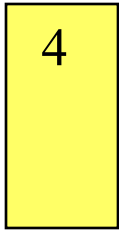


Fill From Disk

Fill From Tree

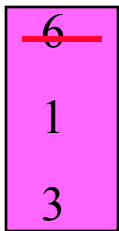
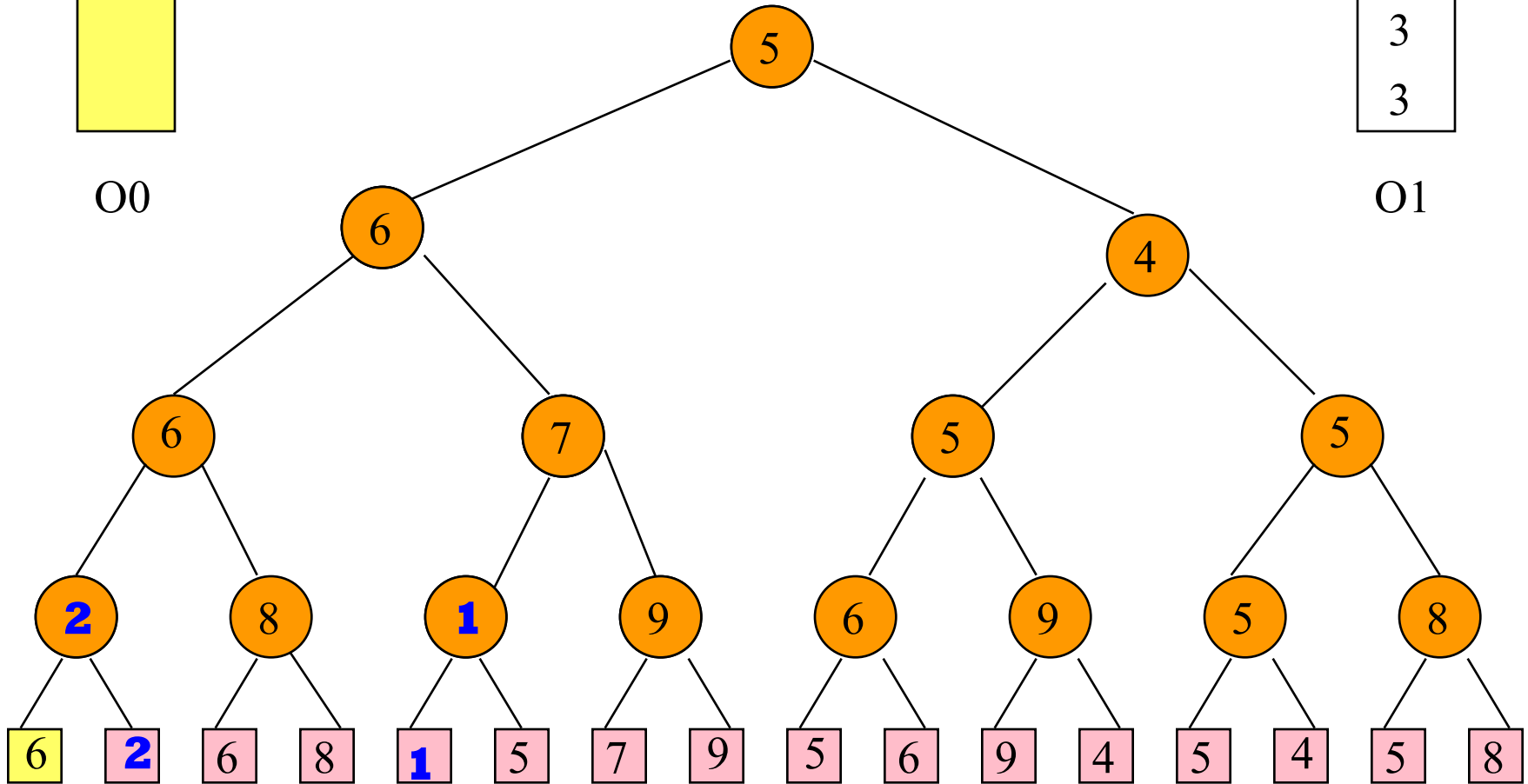
# Continue With Run 1

Write To Disk



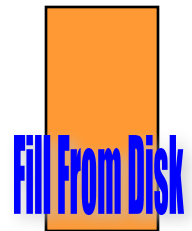
O0

O1



I0

I1

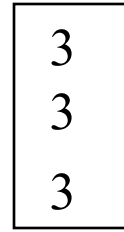
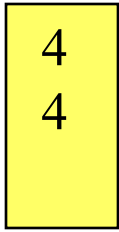


Fill From Disk

Fill From Tree

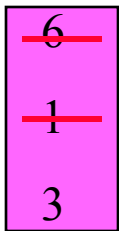
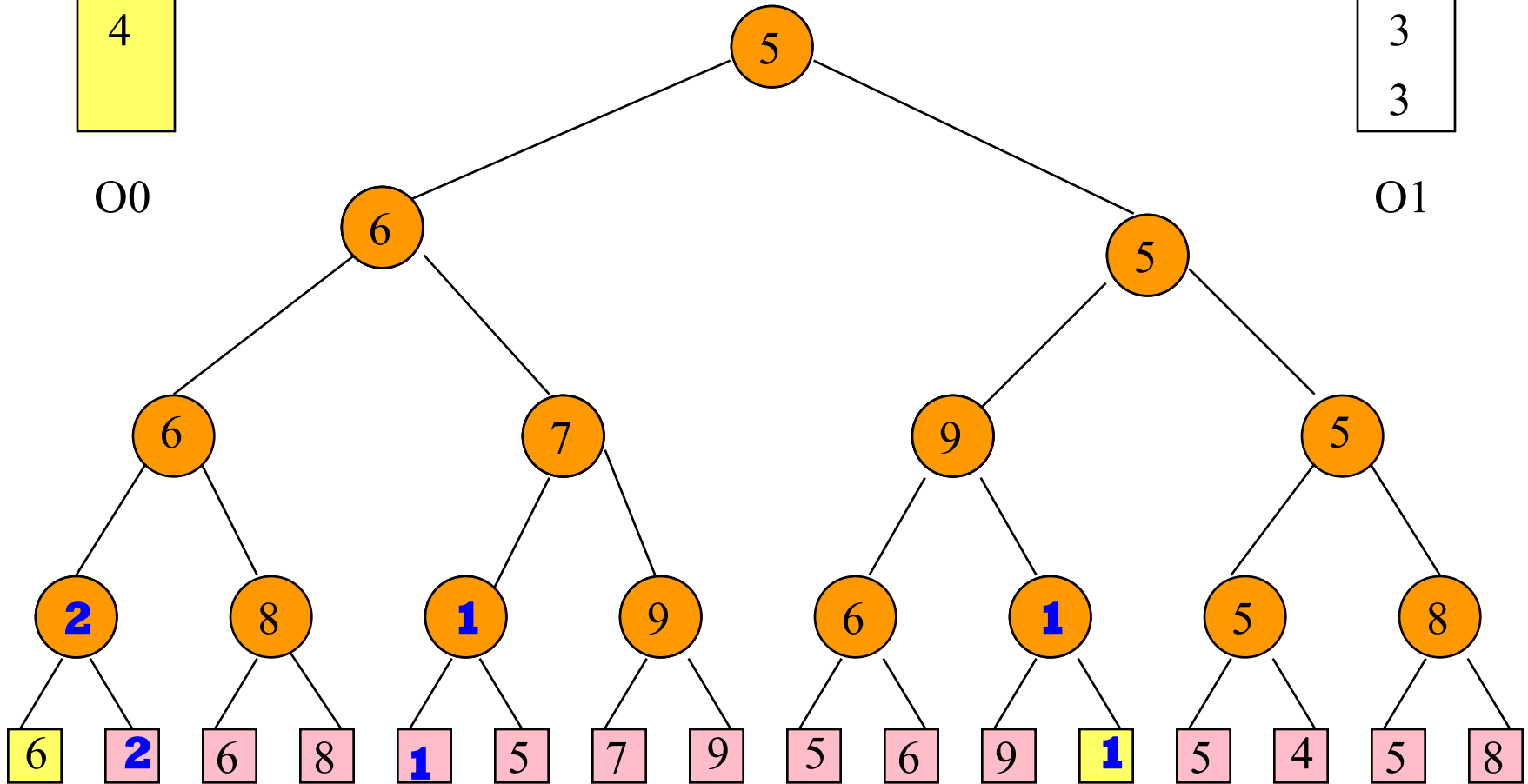
# Continue With Run 1

Write To Disk



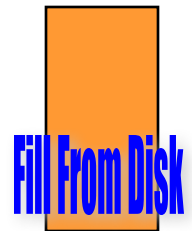
O0

O1



I0

I1



# RUN SIZE

- Let  $k$  be number of external nodes in loser tree.
- Run size  $\geq k$ .
- Sorted input  $\Rightarrow 1$  run.
- Reverse of sorted input  $\Rightarrow n/k$  runs.
- Average run size is  $\sim 2k$ .

# Comparison

- Memory capacity =  $m$  records.
- Run size using fill memory, sort, and output run scheme =  $m$ .
- Use loser tree scheme.
  - Assume block size is  $b$  records.
  - Need memory for 4 buffers ( $4b$  records).
  - Loser tree  $k = m - 4b$ .
  - Average run size =  $2k = 2(m - 4b)$ .
  - $2k \geq m$  when  $m \geq 8b$ .



# Comparison

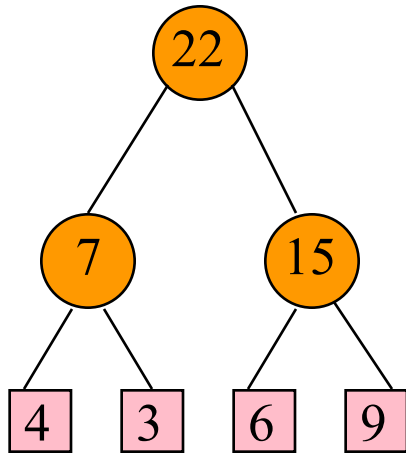
- Assume  $b = 100$ .

$m$	600	1000	5000	10000
$k$	200	600	4600	9600
$2k$	400	1200	9200	19200

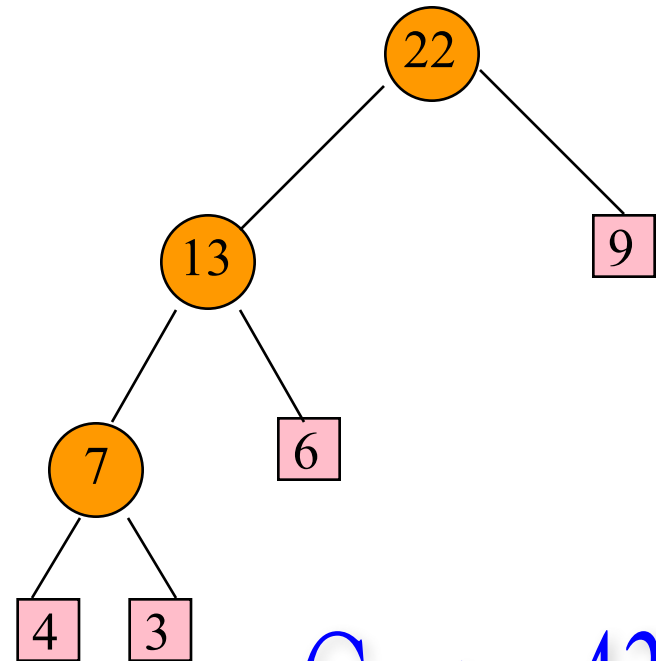
# Comparison

- Total internal processing time using fill memory, sort, and output run scheme  
 $= O((n/m) m \log m) = O(n \log m)$ .
- Total internal processing time using loser tree  $= O(n \log k)$ .
- Loser tree scheme generates runs that differ in their lengths.

# Merging Runs Of Different Length



Cost = 44



Cost = 42

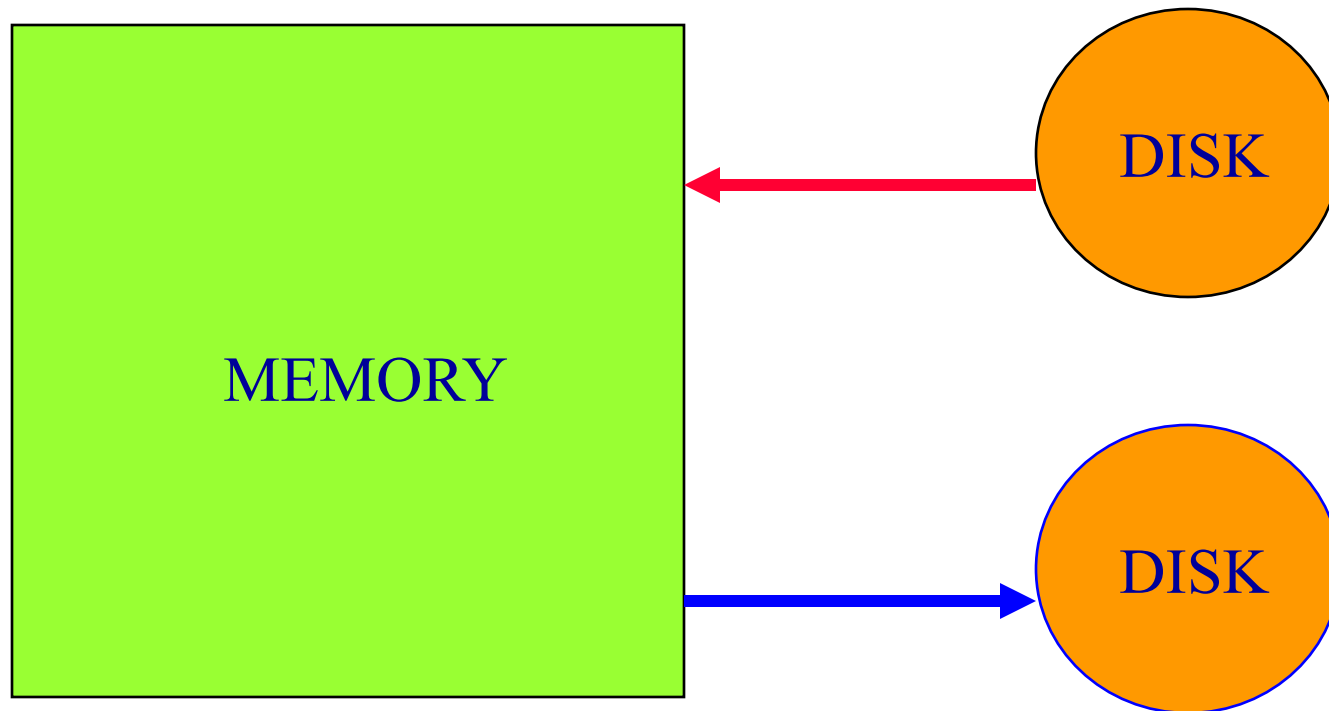
**Best merge sequence?**

# Improve Run Merging

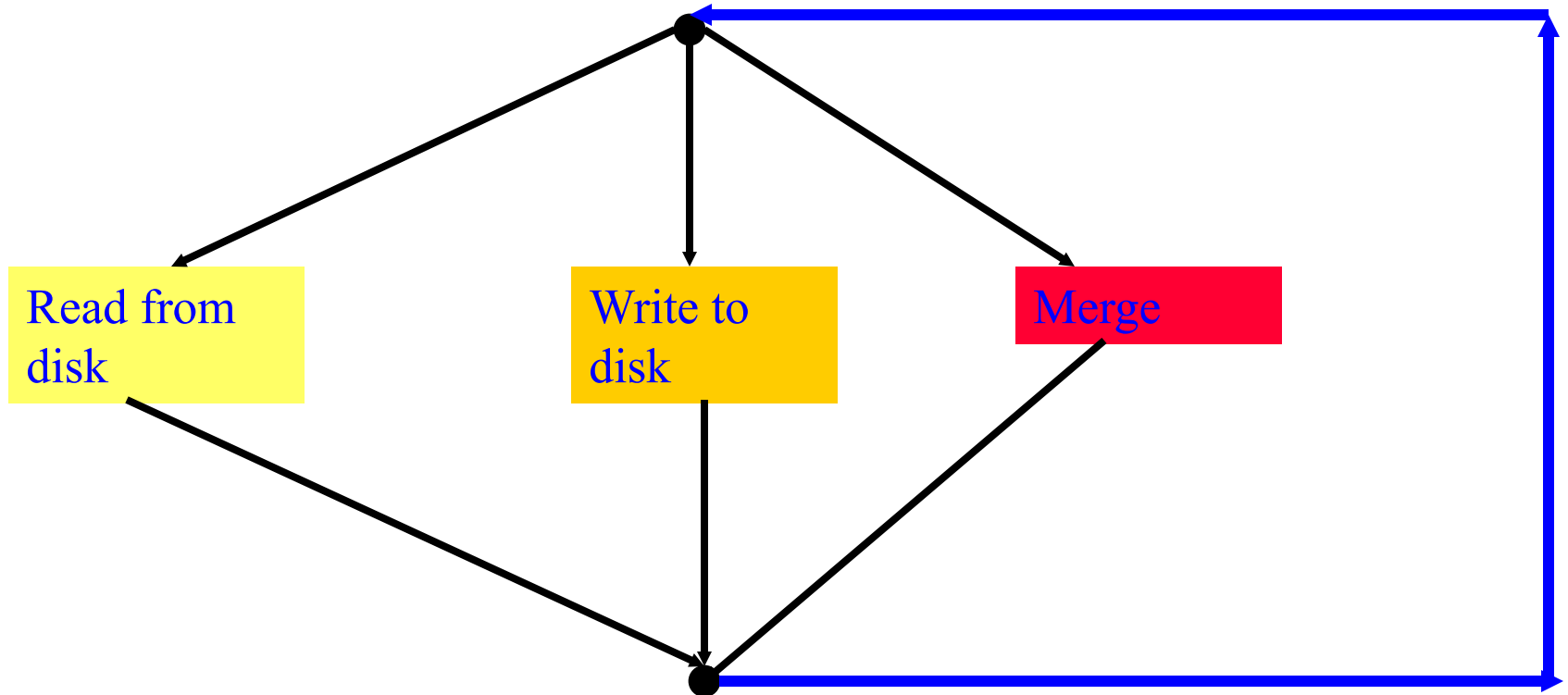
- Reduce number of merge passes.
  - Use higher order merge.
  - Number of passes  
=  $\text{ceil}(\log_k(\text{number of initial runs}))$   
where  $k$  is the merge order.
- More generally, a higher-order merge reduces the cost of the optimal merge tree.

# Improve Run Merging

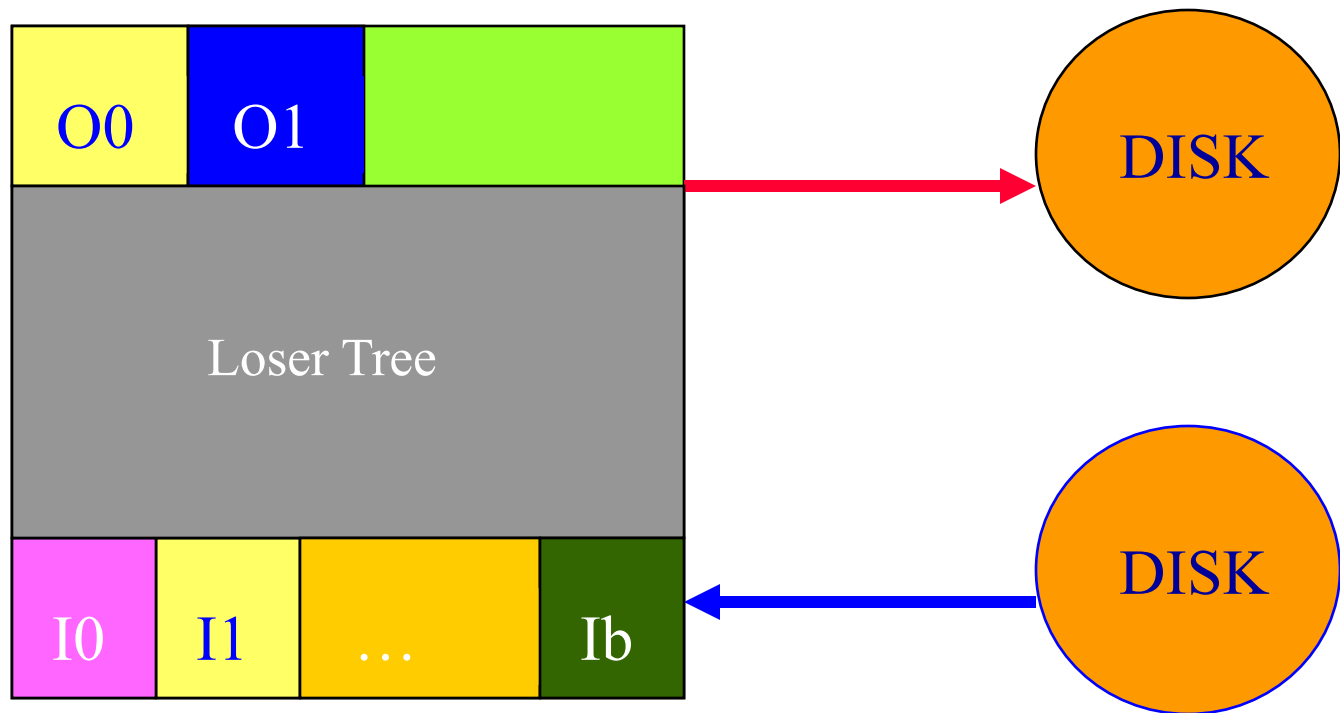
- Overlap input, output, and internal merging.



# Steady State Operation



# Partitioning Of Memory



- Need exactly **2** output buffers.
- Need at least  **$k+1$**  ( **$k$**  is merge order) input buffers.
- **$2k$**  input buffers suffice.

# Number Of Input Buffers

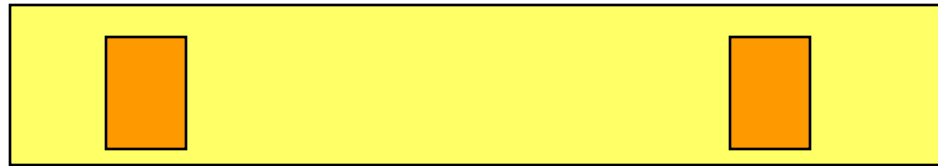
- When 2 input buffers are dedicated to each of the  $k$  runs being merged,  $2k$  buffers are not enough!
- Input buffers must be allocated to runs on an as needed basis.



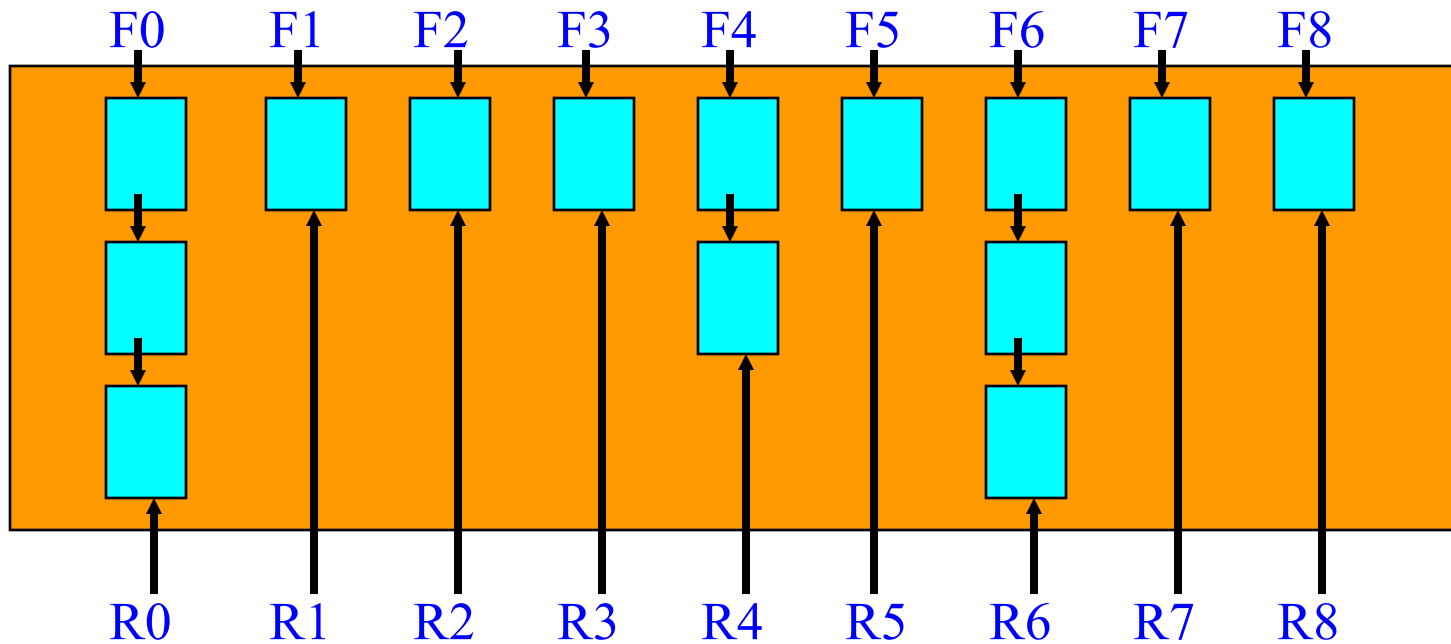
# Buffer Allocation

- When ready to read a buffer load, determine which run will exhaust first.
  - Examine key of the last record read from each of the  $k$  runs.
  - Run with smallest last key read will exhaust first.
- Next buffer load of input is to come from run that will exhaust first, allocate an input buffer to this run.

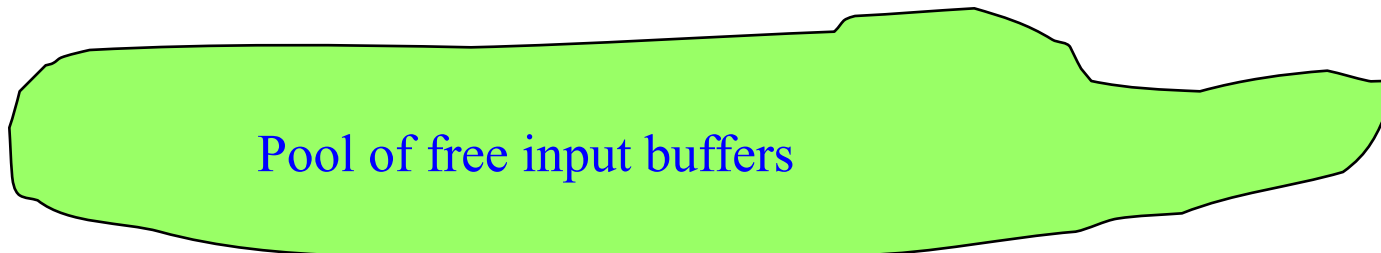
# Buffer Layout



Output  
buffers



Input buffer  
queues  $k=9$



Pool of free input buffers

# Initialize To Merge $k$ Runs

- Initialize  $k$  queues of input buffers, 1 queue per run, 1 buffer per run.
- Input one buffer load from each of the  $k$  runs.
- Put  $k - 1$  unused input buffers into pool of free buffers.
- Set  $\text{activeOutputBuffer} = 0$ .
- Initiate input of next buffer load from first run to exhaust. Use remaining unused input buffer for this input.

# The Method kWayMerge

- k-way merge from input queues to the active output buffer.
- Merge stops when either the output buffer gets full or when an end-of-run key is merged into the output buffer.
- If merge hasn't stopped and an input buffer gets empty, advance to next buffer in queue and free empty buffer.

# Merge k Runs

repeat

kWayMerge;

wait for input/output to complete;

add new input buffer (if any) to queue for its run;

determine run that will exhaust first;

if (there is more input from this run)

initiate read of next block for this run;

initiate write of active output buffer;

activeOutputBuffer = 1 – activeOutputBuffer;

until end-of-run key merged;

# What Can Go Wrong?

## kWayMerge

- k-way merge from input queues to the active output buffer.
- Merge stops when either the output buffer gets full or when an end-of-run key is merged into the output buffer.
- If merge hasn't stopped and an input buffer gets empty, advance to next buffer in queue and free empty buffer.  
There may be no next buffer in the queue.

# What Can Go Wrong?

## Merge k Runs

repeat

kWayMerge;

wait for input/output to complete;

add new input buffer (if any) to queue for its run;

determine run that will exhaust first;

if (there is more input from this run)

initiate read of next block for this run;

There may be no  
free input buffer  
to read into.

initiate write of active output buffer;

$\text{activeOutputBuffer} = 1 - \text{activeOutputBuffer};$

until end of run key merged;

# kWayMerge

- If merge hasn't stopped and an input buffer gets empty, advance to next buffer in queue and free empty buffer.  
There may be no next buffer in the queue.

- If this type of failure were to happen, using two different and valid analyses, we will end up with inconsistent counts of the amount of data available to kWayMerge.
- Data available to kWayMerge is data in
  - Input buffer queues.
  - Active output buffer.
  - Excludes data in buffer being read or written.



# No Next Buffer In Queue

repeat

kWayMerge; 

wait for input/output to complete;

add new input buffer (if any) to queue for its run;

determine run that will exhaust first;

if (there is more input from this run)

initiate read of next block for this run;

initiate write of active output buffer;

activeOutputBuffer = 1 – activeOutputBuffer;

until end-of-run key merged;

- Exactly **k** buffer loads available to kWayMerge.

# kWayMerge

- If merge hasn't stopped and an input buffer gets empty, advance to next buffer in queue and free empty buffer.  
There may be no next buffer in the queue.
- Alternative analysis of data available to kWayMerge at time of failure.
  - $< 1$  buffer load in active output buffer
  - $\leq k - 1$  buffer loads in remaining  $k - 1$  queues
  - Total data available to k-way merge is  $< k$  buffer loads.

# Merge $k$ Runs

initiate read of next block for this run;

There may be no free input buffer to read into.

- Suppose there is no free input buffer.
- One analysis will show there are exactly  $k + 1$  buffer loads in memory (including newly read input buffer) at time of failure.
- Another analysis will show there are  $> k + 1$  buffer loads in memory at time of failure.
- Note that at time of failure there is no buffer being read or written.

# No Free Input Buffer

repeat

kWayMerge;

wait for input/output to complete;

add new input buffer (if any) to queue for its run;

determine run that will exhaust first;

if (there is more input from this run)

initiate read of next block for this run;



initiate write of active output buffer;

activeOutputBuffer = 1 – activeOutputBuffer;

until end-of-run key merged;

- Exactly  $k + 1$  buffer loads in memory.

# Merge $k$ Runs

initiate read of next block for this run;

There may be no free input buffer to read into.

- Alternative analysis of data in memory.
  - 1 buffer load in the active output buffer.
  - 1 input queue may have an empty first buffer.
  - Remaining  $k - 1$  input queues have a nonempty first buffer.
  - Remaining  $k$  input buffers must be in queues and full.
  - Since  $k > 1$ , total data in memory is  $> k + 1$  buffer loads.

# Minimize Wait Time For I/O To Complete

Time to fill an output buffer

~ time to read a buffer load

~ time to write a buffer load

# Initializing For Next k-way Merge

Change

**if** (there is more input from this run)

    initiate read of next block for this run;

to

**if** (there is more input from this run)

    initiate read of next block for this run;

**else**

    initiate read of a block for the next k-way merge;

# What Can Go Wrong?

## kWayMerge

- k-way merge from input queues to the active output buffer.
- Merge stops when either the output buffer gets full or when an end-of-run key is merged into the output buffer.
- If merge hasn't stopped and an input buffer gets empty, advance to next buffer in queue and free empty buffer.  
There may be no next buffer in the queue.



# What Can Go Wrong?

## Merge k Runs

repeat

kWayMerge;

wait for input/output to complete;

add new input buffer (if any) to queue for its run;

determine run that will exhaust first;

if (there is more input from this run)

initiate read of next block for this run;

There may be no  
free input buffer  
to read into.

initiate write of active output buffer;

$\text{activeOutputBuffer} = 1 - \text{activeOutputBuffer};$

until end of run key merged;

# Initializing For Next k-way Merge

Change

**if** (there is more input from this run)

    initiate read of next block for this run;

to

**if** (there is more input from this run)

    initiate read of next block for this run;

**else**

    initiate read of a block for the next k-way merge;