

Today

Distributed Subsystems (Lesson 7)

- ⇒ \* Global memory System
- ⇒ \* Distributed Shared Memory

Wednesday

\* Distributed file Systems

## Question

Fill in the table with what happens to the boundary between Local and Global on each node

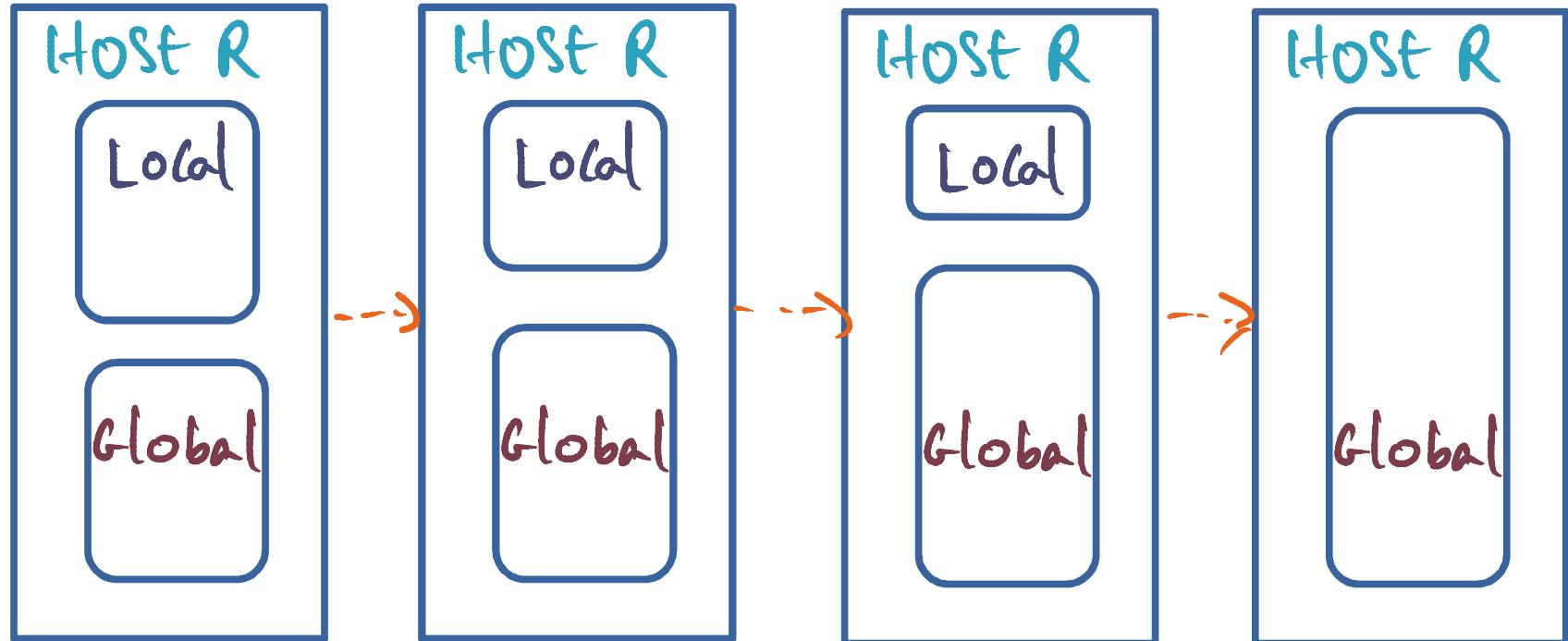
Faulting Page X	Faulting node P L G	Node Q with Page X L G	Node R with LRU page L G
in Q's global			
in Q's global P's global empty			
on disk			
actively shared with Q			

## Question

Fill in the table with what happens to the boundary between Local and Global on each node

Faulting Page X	Faulting node P L G	Node Q with page X L G	Node R with LRU page L G
in Q's global	+1 -1	no change	no change
in Q's global P's global empty	no change	no change	no change
on disk	+1 -1	Not applicable	<u>-1 +1</u>
actively shared with Q	+1 -1	no change	<u>-1 +1</u>

## Behavior of Algorithm



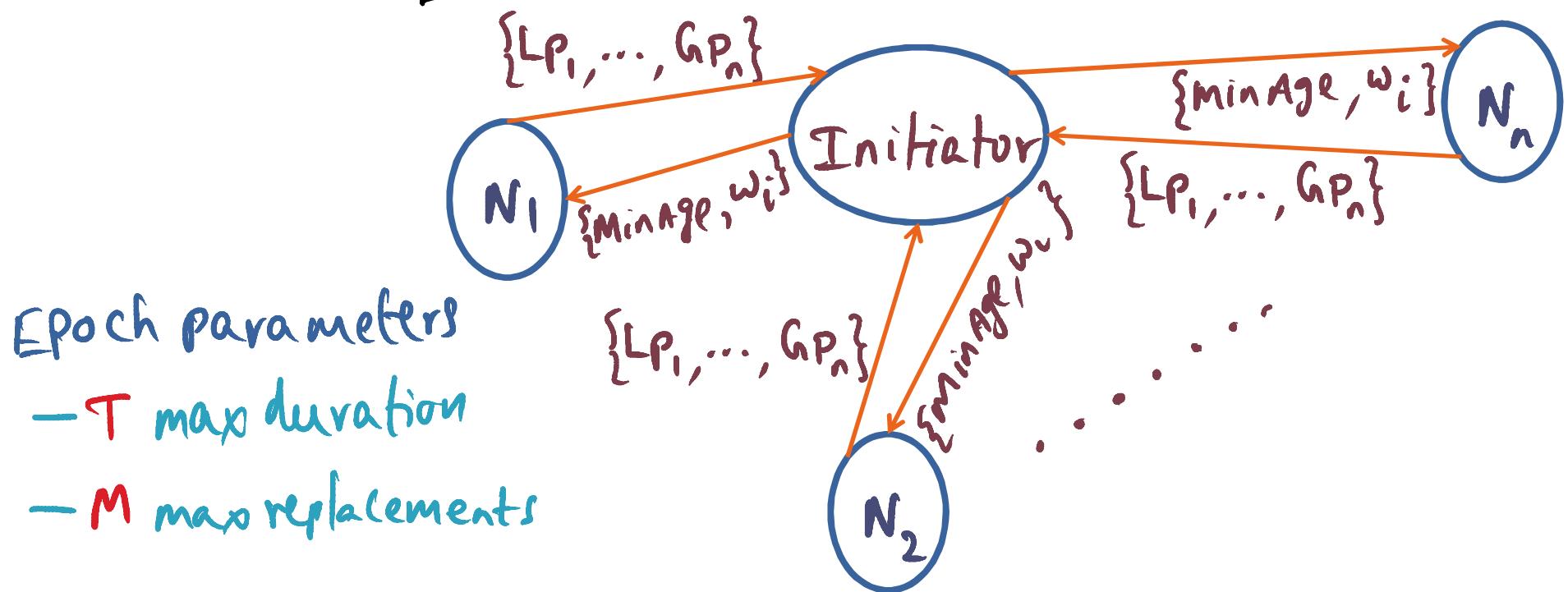
R becomes a memory server for peers  
on the cluster

# Geriatrics!!

Epoch parameters

- $T$  max duration
- $M$  max replacements

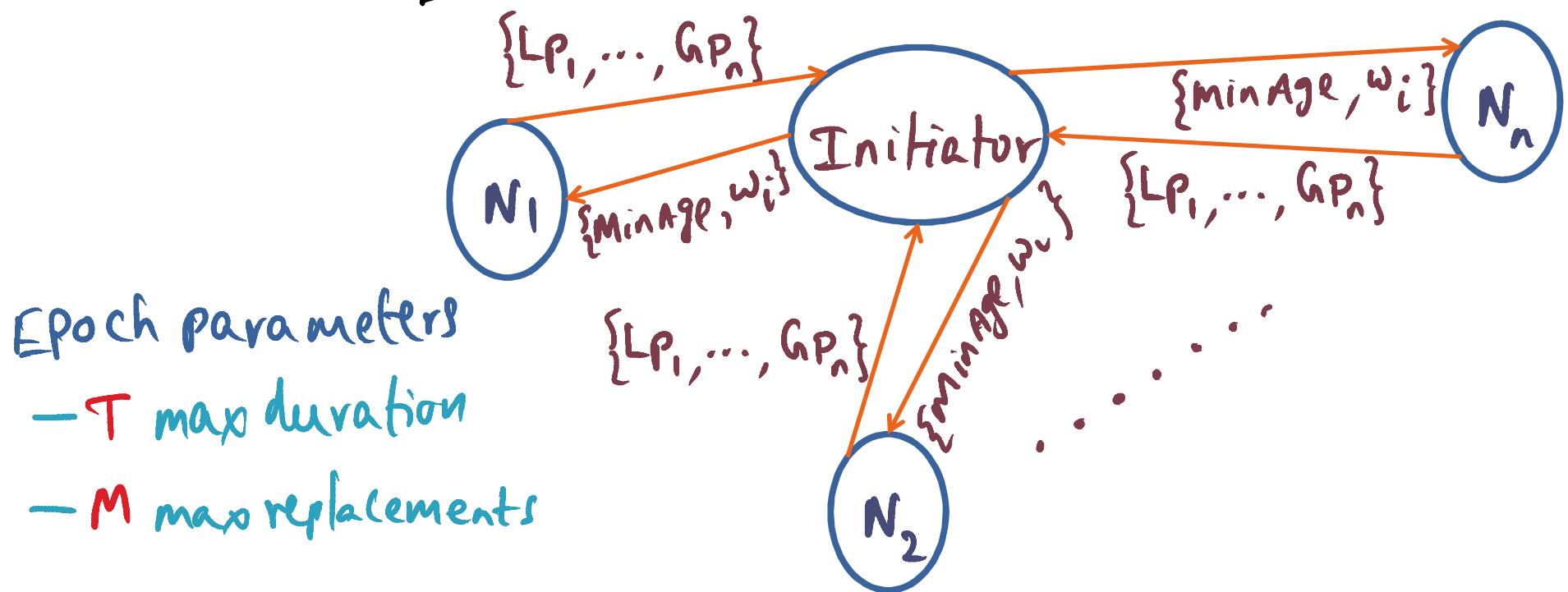
## Geriatrics !!



Each Epoch

- Send age info to initiator
- receive  $\{\text{minAge}, w_i\}$  + i

## Geriatrics !!



Epoch parameters

- $T$  max duration
- $M$  max replacements

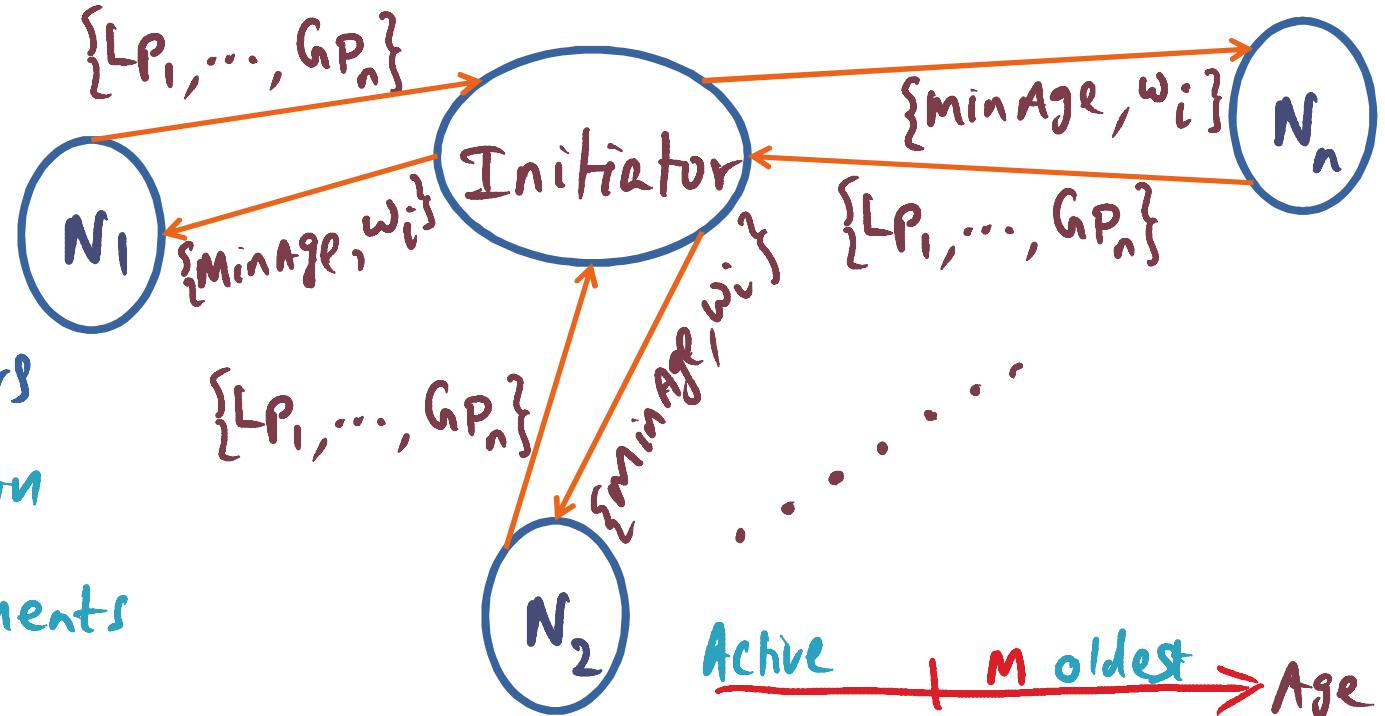
Each Epoch

- Send age info to initiator
- receive  $\{minAge, w_i\}$  +  $i$

Initiator for next epoch

- node with  $\max(w_i)$

## Geriatrics !!



Epoch parameters

- $T$  max duration
- $M$  max replacements

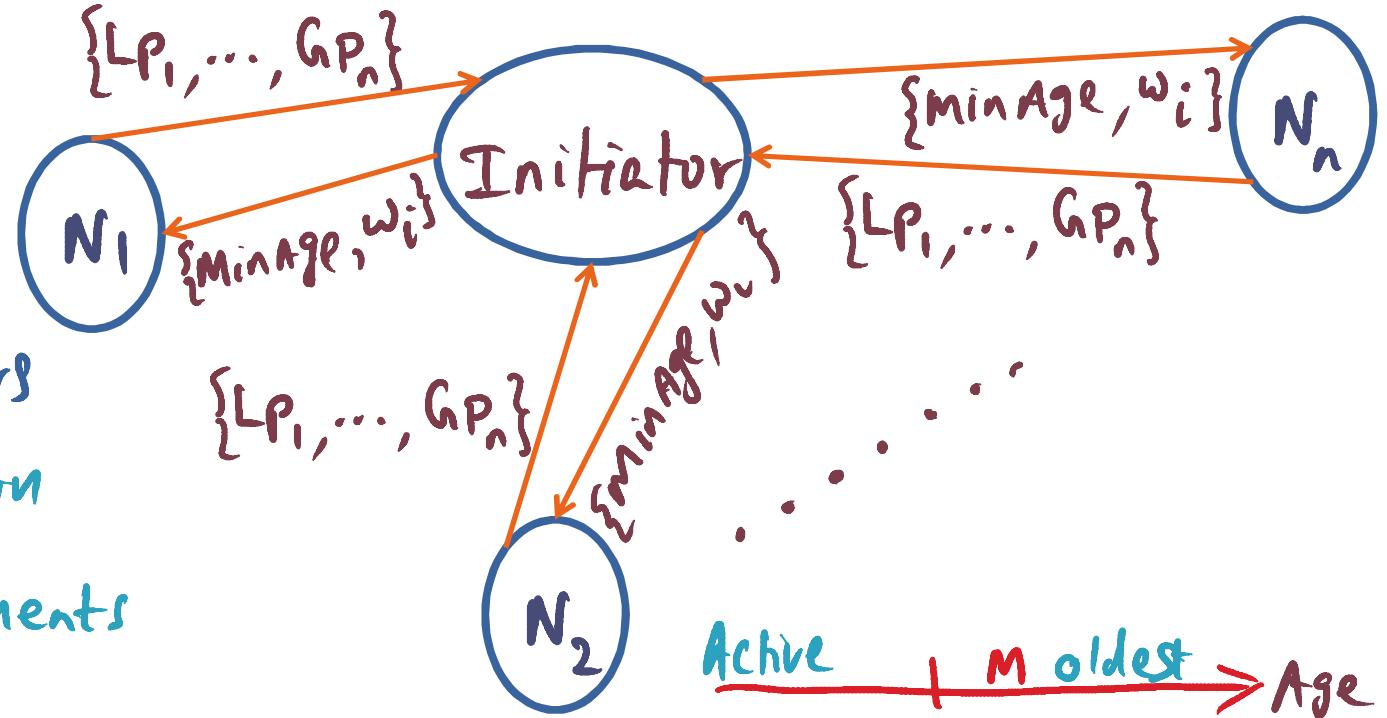
Each Epoch

- Send age info to initiator
- receive  $\{\text{MinAge}, w_i\}$

Initiator for next epoch

- node with  $\text{Max}(w_i)$

# Geriatrics !!



Epoch parameters

- $T$  max duration
- $M$  max replacements

Each Epoch

- Send age info to initiator
- receive  $\{\text{MinAge}, w_i\}$

Initiator for next epoch

- node with  $\text{Max}(w_i)$

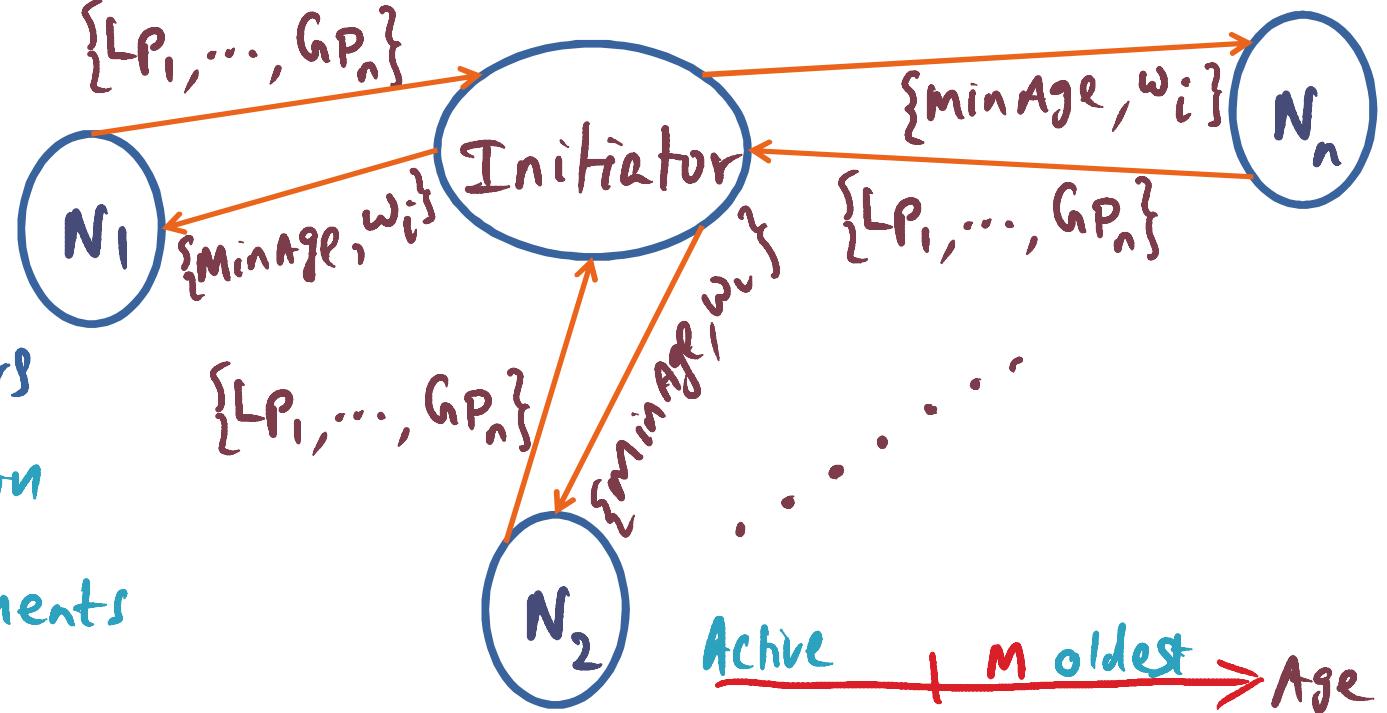
Action at a node on page fault  
Page y eviction candidate

$\text{Age}(\text{page } y) > \text{MinAge} \Rightarrow \text{discard}$

$\text{Age}(\text{page } y) < \text{MinAge}$   
 $\Rightarrow \text{Send to peer } N_i$

Think Global  
Act local !!

Geriatrics !!



Epoch parameters

- $T$  max duration
- $M$  max replacements

Each Epoch

- Send age info to initiator
- receive  $\{\text{MinAge}, w_i\}$

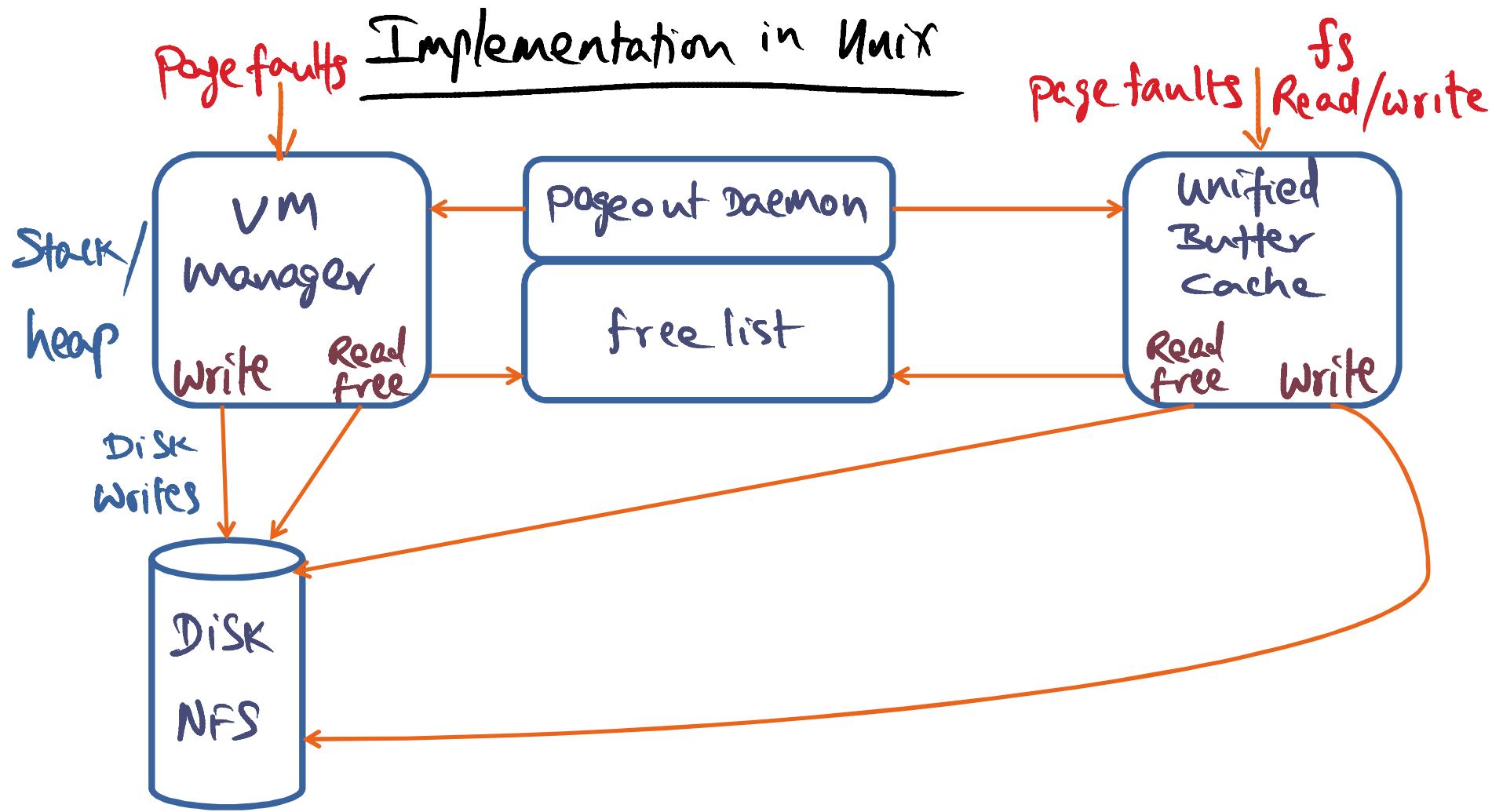
Initiator for next epoch

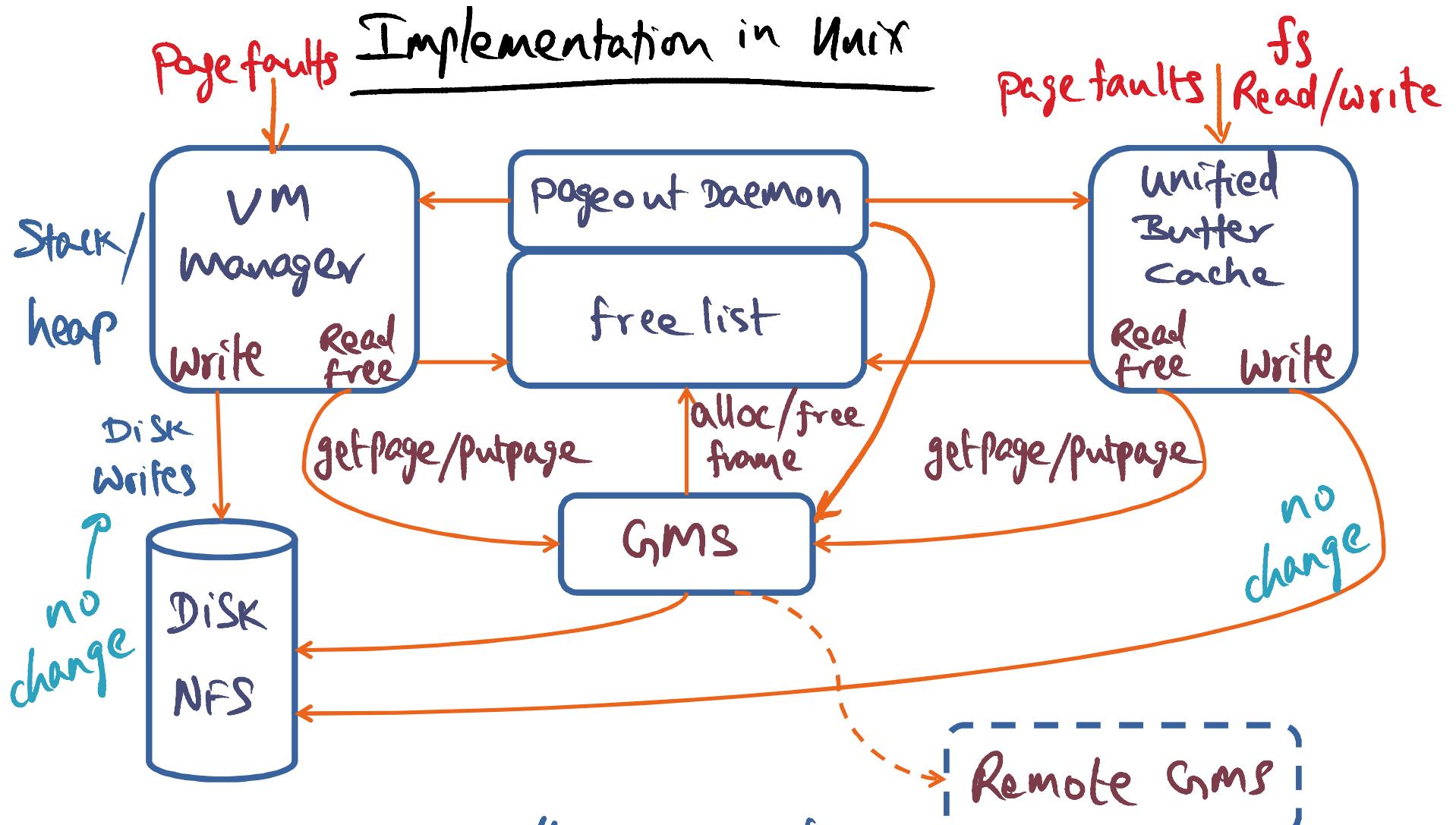
- node with  $\text{Max}(w_i)$

Action at a node on page fault  
Page y eviction candidate

$\text{Age}(\text{page } y) > \text{MinAge} \Rightarrow \text{discard}$

$\text{Age}(\text{page } y) < \text{MinAge}$   
 $\Rightarrow \text{Send to peer } N_i$





GMS integrated with DEC OSF/i

- access to anonymous pages + f.s. mapped

Pages go through GMS on reads

## Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

Data structure for

UID  $\rightarrow$  physical page frame ?

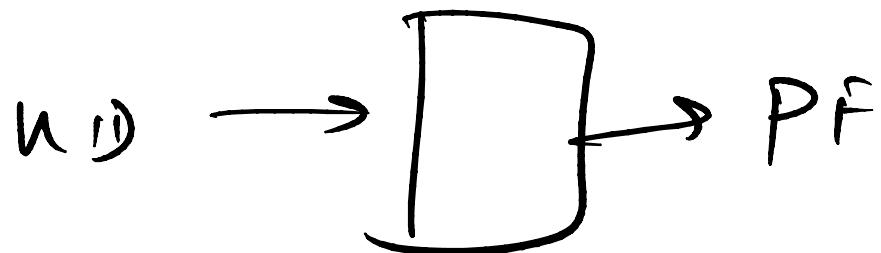
# Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

Data structure for  
UID  $\rightarrow$  physical page frame ?



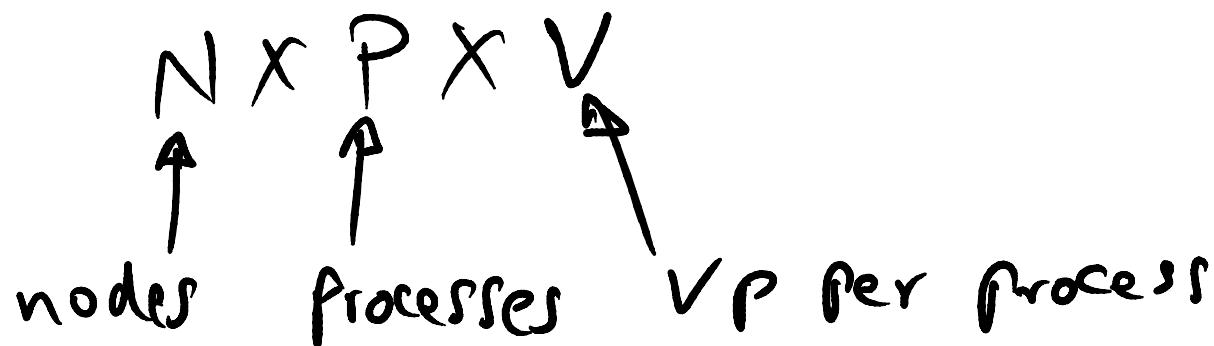
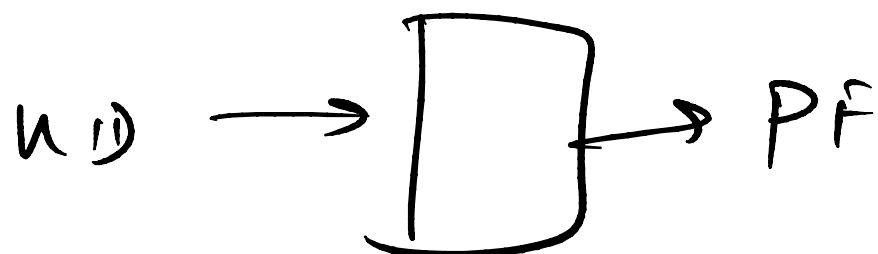
# Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

Data structure for  
UID  $\rightarrow$  physical page frame ?



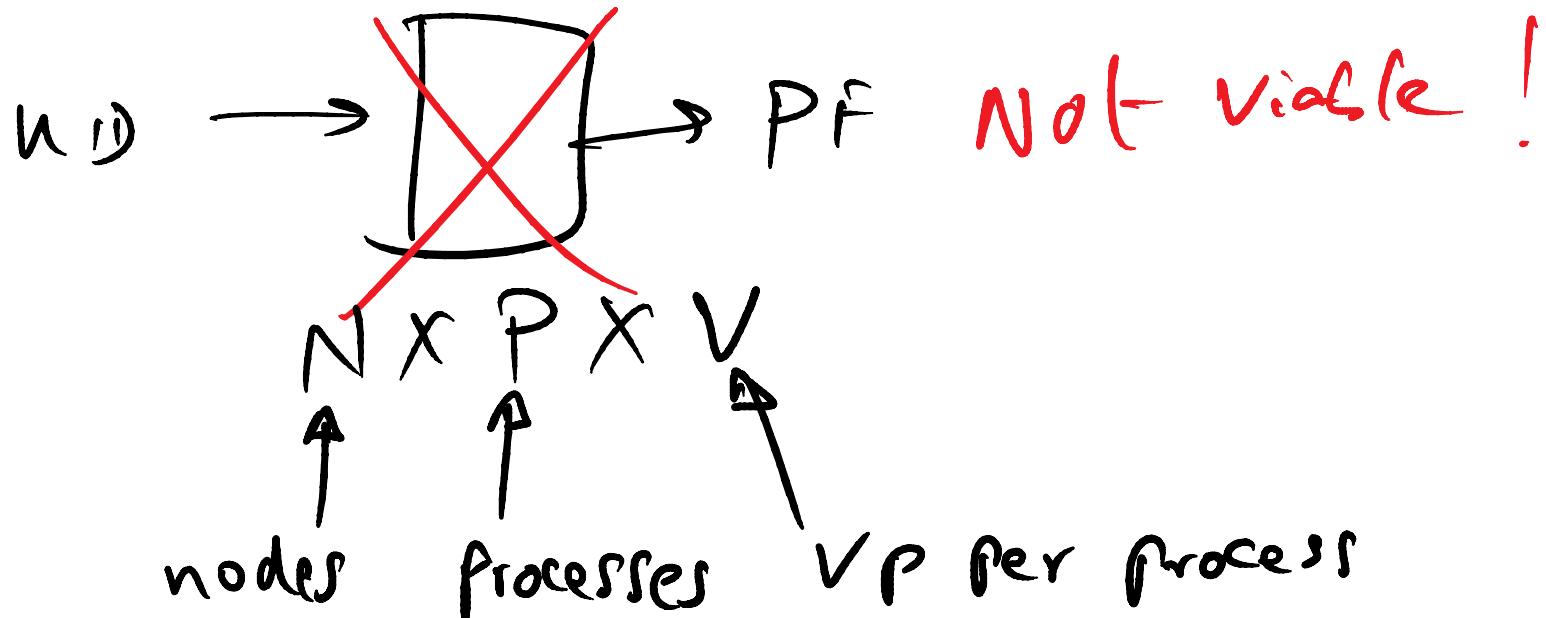
# Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

Data structure for  
UID  $\rightarrow$  physical page frame ?



## Twin Goals

- \* reduce communication
- \* scalable with number of nodes
  - ⇒ don't burden any one node

## Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

Data structure for

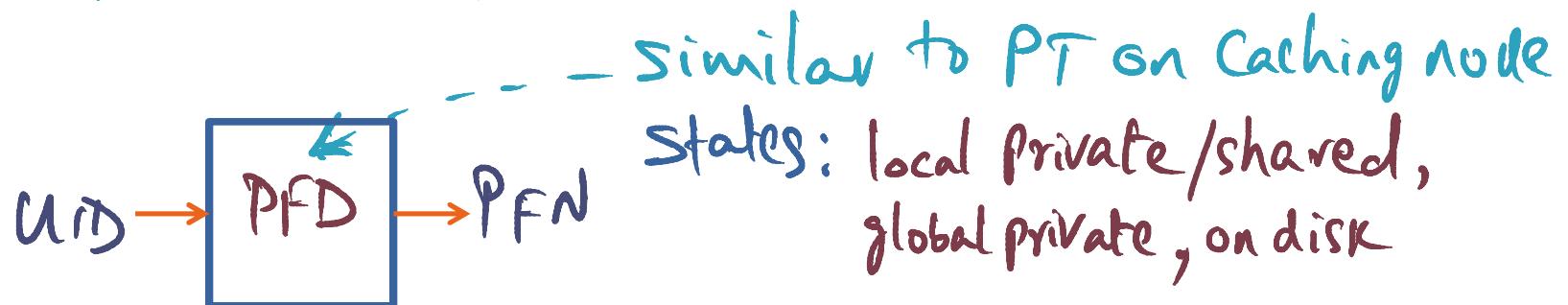
UID  $\rightarrow$  physical page frame ?

# Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

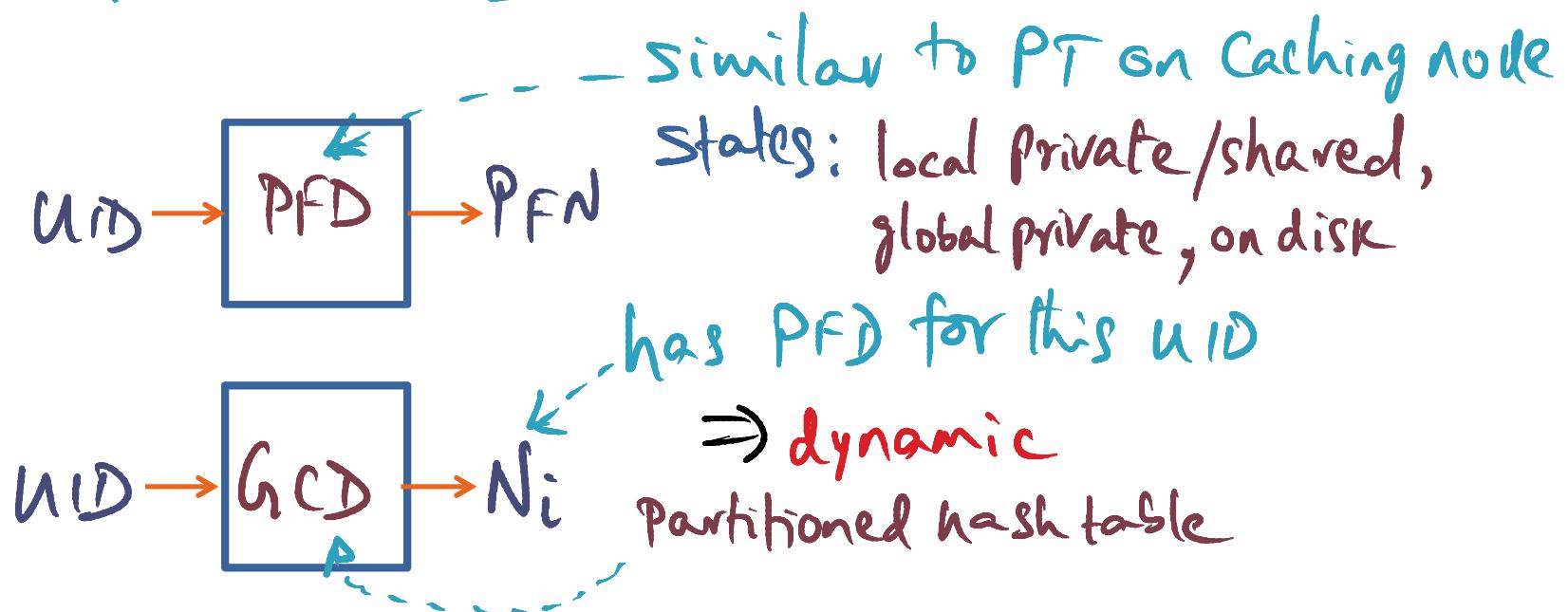


# Data structures

VA  $\rightarrow$  UID

IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC

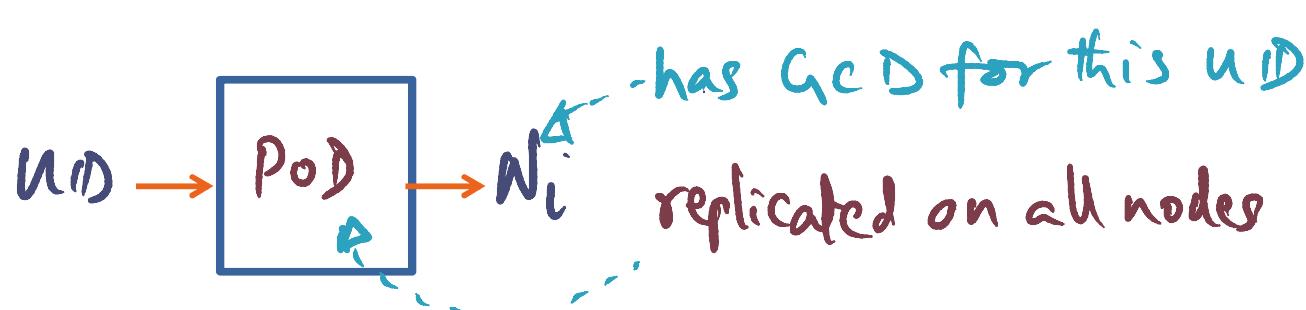
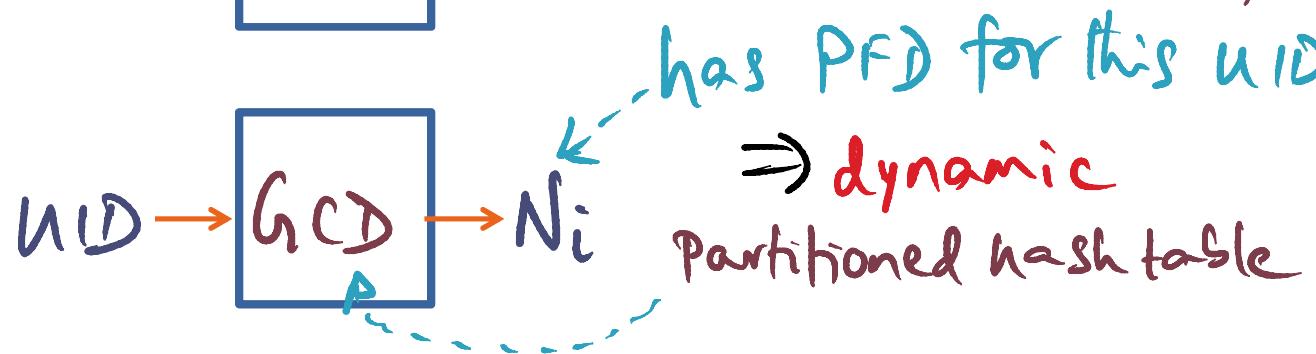
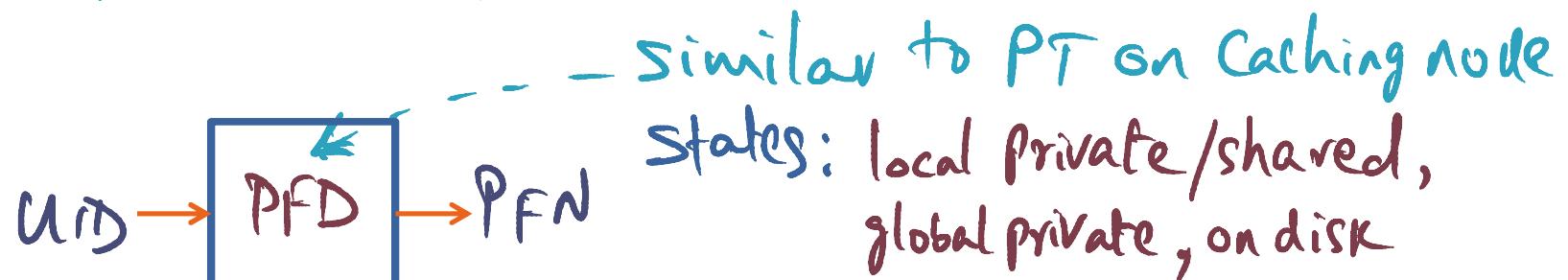


# Data structures

VA  $\rightarrow$  UID

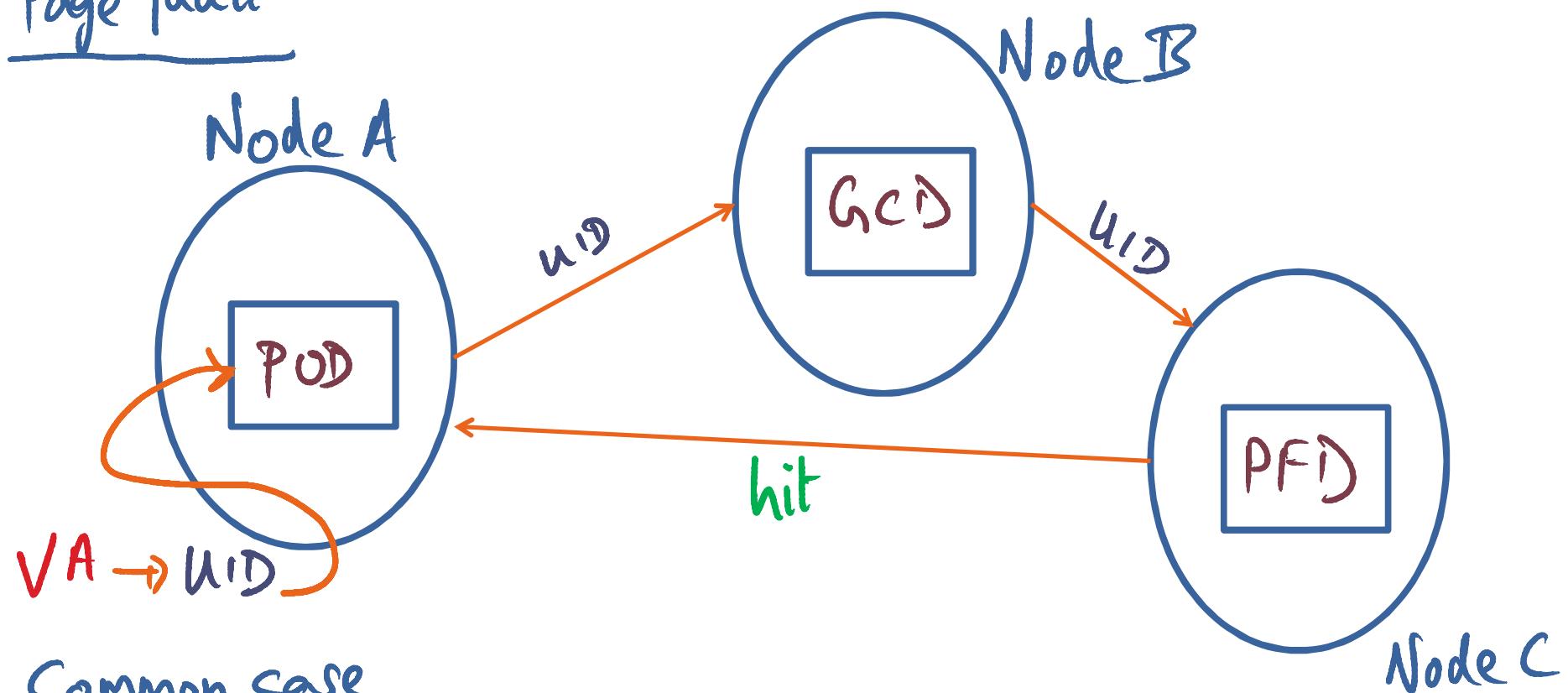
IP-Addr	disk partition	i-node	offset
---------	----------------	--------	--------

- derived from VM + UBC



# Putting the Data structures to work

Page fault

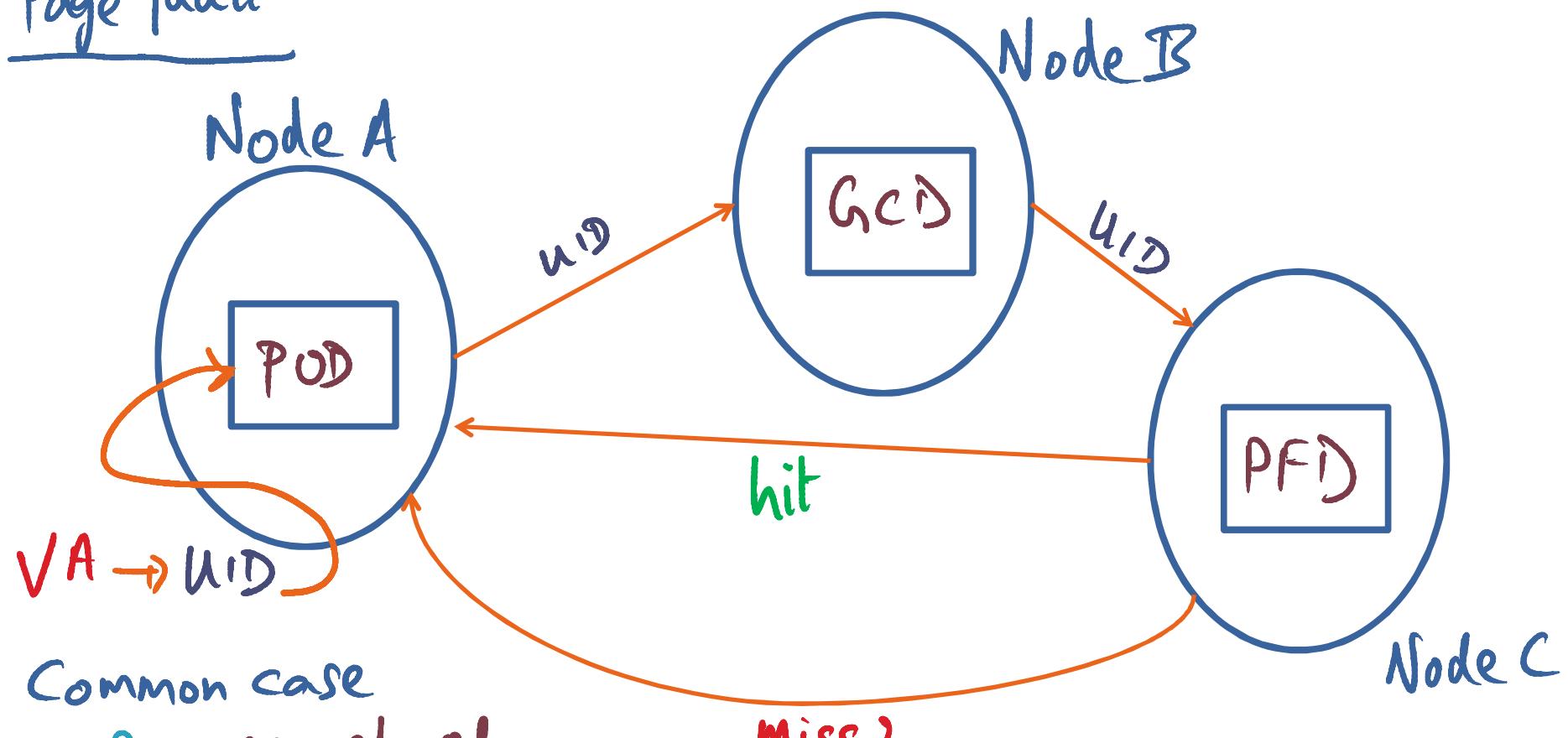


Common case

- Page non-shared  
 $\Rightarrow A + B$  same
- Page fault service quick

# Putting the Data structures to work

Page fault

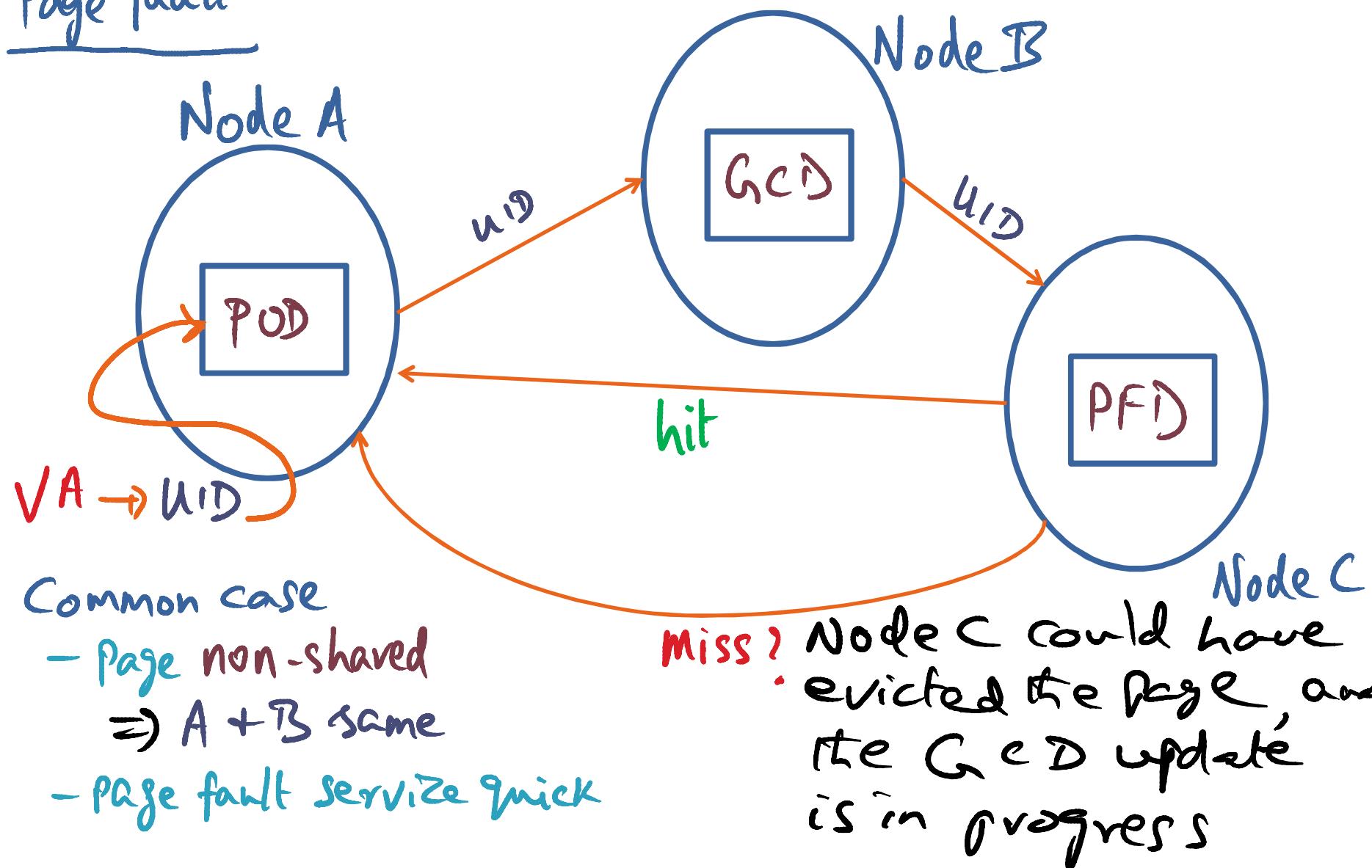


Common case

- Page non-shared  
 $\Rightarrow A + B$  same
- Page fault service quick

# Putting the Data structures to work

Page fault



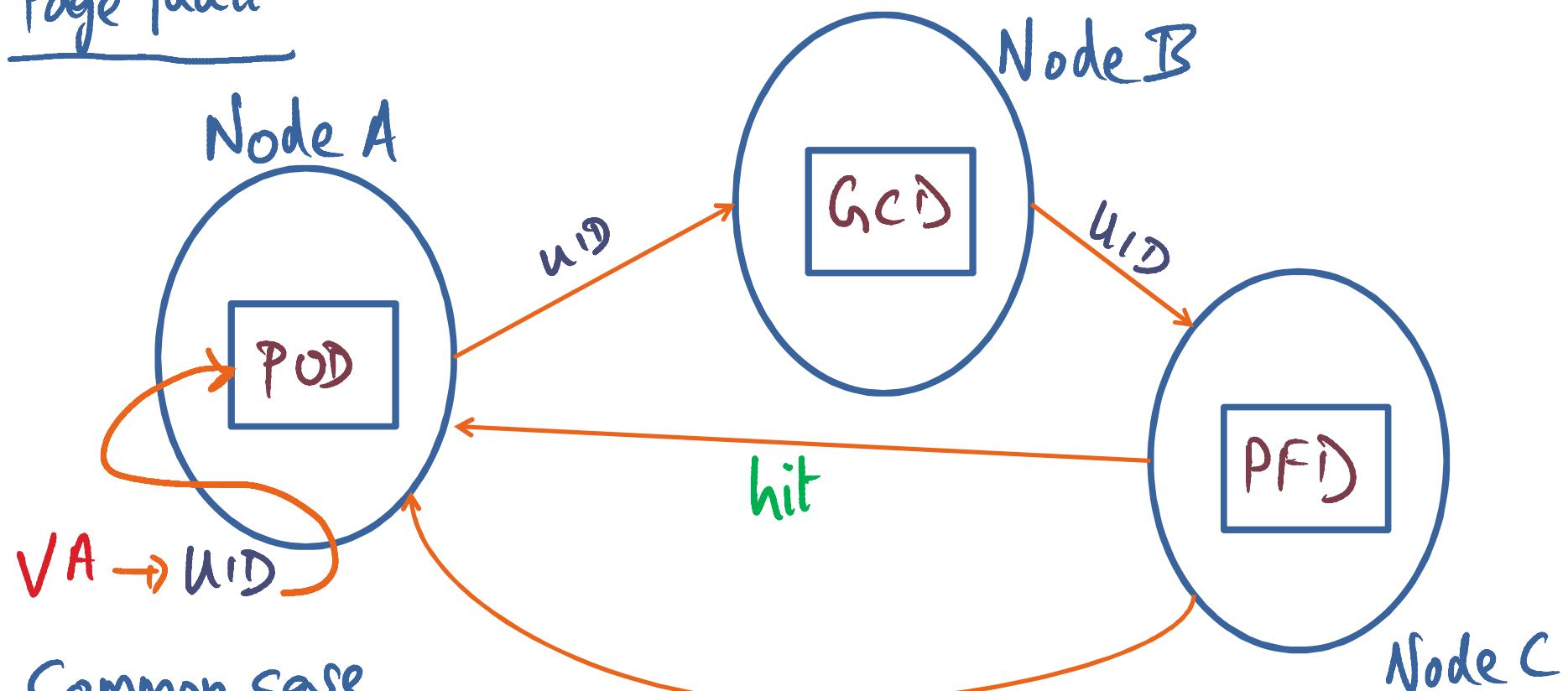
Common case

- Page non-shared  
 $\Rightarrow A + B$  same
- Page fault service quick

Miss? Node C could have evicted the page, and the GCD update is in progress

# Putting the Data structures to work

Page fault



Common case

- Page non-shared
- $\Rightarrow A + B$  same
- Page fault service quick

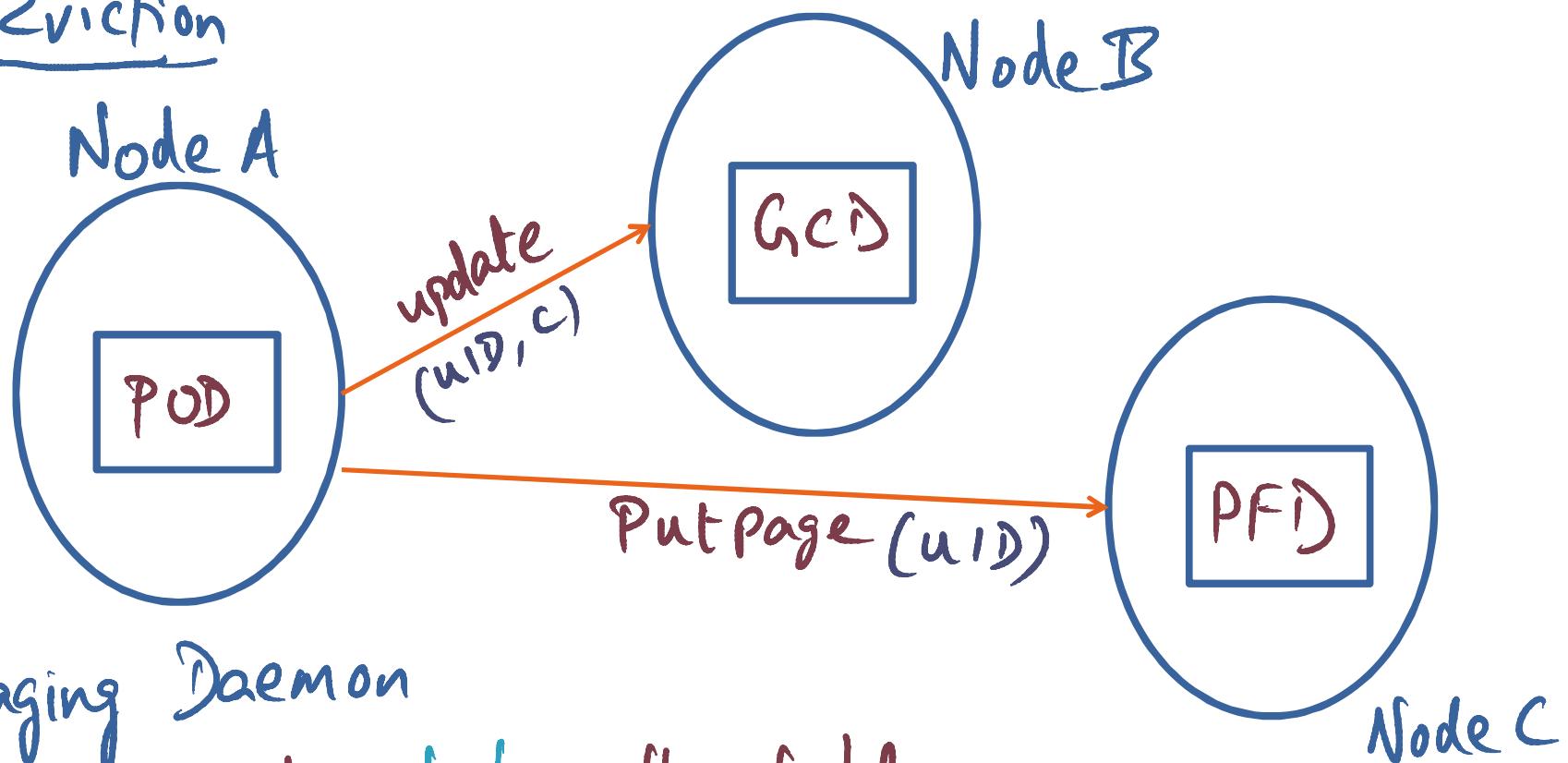
Miss?

$\Rightarrow$  Uncommon Case

- POD changing due to node addition/deletion

# Putting the Data structures to work

## Page Eviction



## Paging Daemon

- Freelist below threshold
- Put page oldest pages
- update GCD, PFD for the UIDs

## Question

Fill in the table with what happens to the boundary between Local and Global on each node

Faulting Page X	Faulting node P L G	Node Q with Page X L G	Node R with LRU page L G
in Q's global			
in Q's global P's global empty			
on disk			
actively shared with Q			

## Question

Fill in the table with what happens to the boundary between Local and Global on each node

Faulting Page X	Faulting node P L G	Node Q with page X L G	Node R with LRU page L G
in Q's global	+1 -1	no change	no change
in Q's global P's global empty	no change	no change	no change
on disk	+1 -1	Not applicable	<u>-1 +1</u>
actively shared with Q	+1 -1	no change	<u>-1 +1</u>

Write down key takeaways and  
give to your neighbor

# Key takeaways

- Heavy lifting to be done to take a **concept** to **implementation**.
  - non-trivial intellectual exercise in building distributed subsystems
- What is **enduring** in a research exercise like this one?
  - The concept of paging across the network is an interesting thought experiment but it may not be feasible exactly for the environment in which the authors carried out the research (workstation clusters connected to a LAN). Each workstation in that setting is a private resource of an individual who may or may not want his resources (memory in particular) shared with others...
  - On the other hand, **data centers** are powered by large-scale clusters (> 1000s of processors in a LAN). No node is individually “owned” by a user. Could this idea of **paging across the LAN** be feasible in that setting? Perhaps.
- Even beyond the thought experiment itself, what is perhaps more enduring are the techniques (**distributed data structures** and **algorithms**) for taking the concept to implementation