# Open Discussion
# File Systems

Renaud Lottiaux
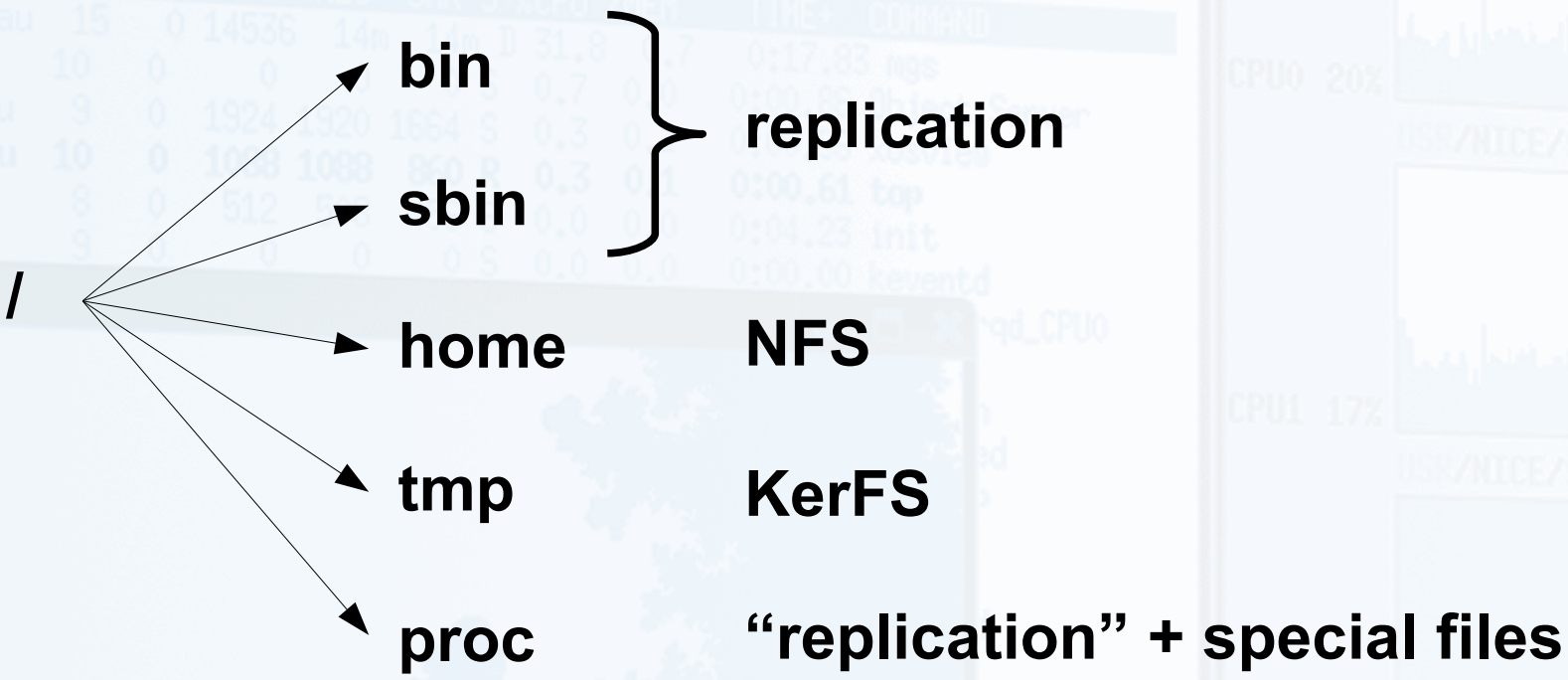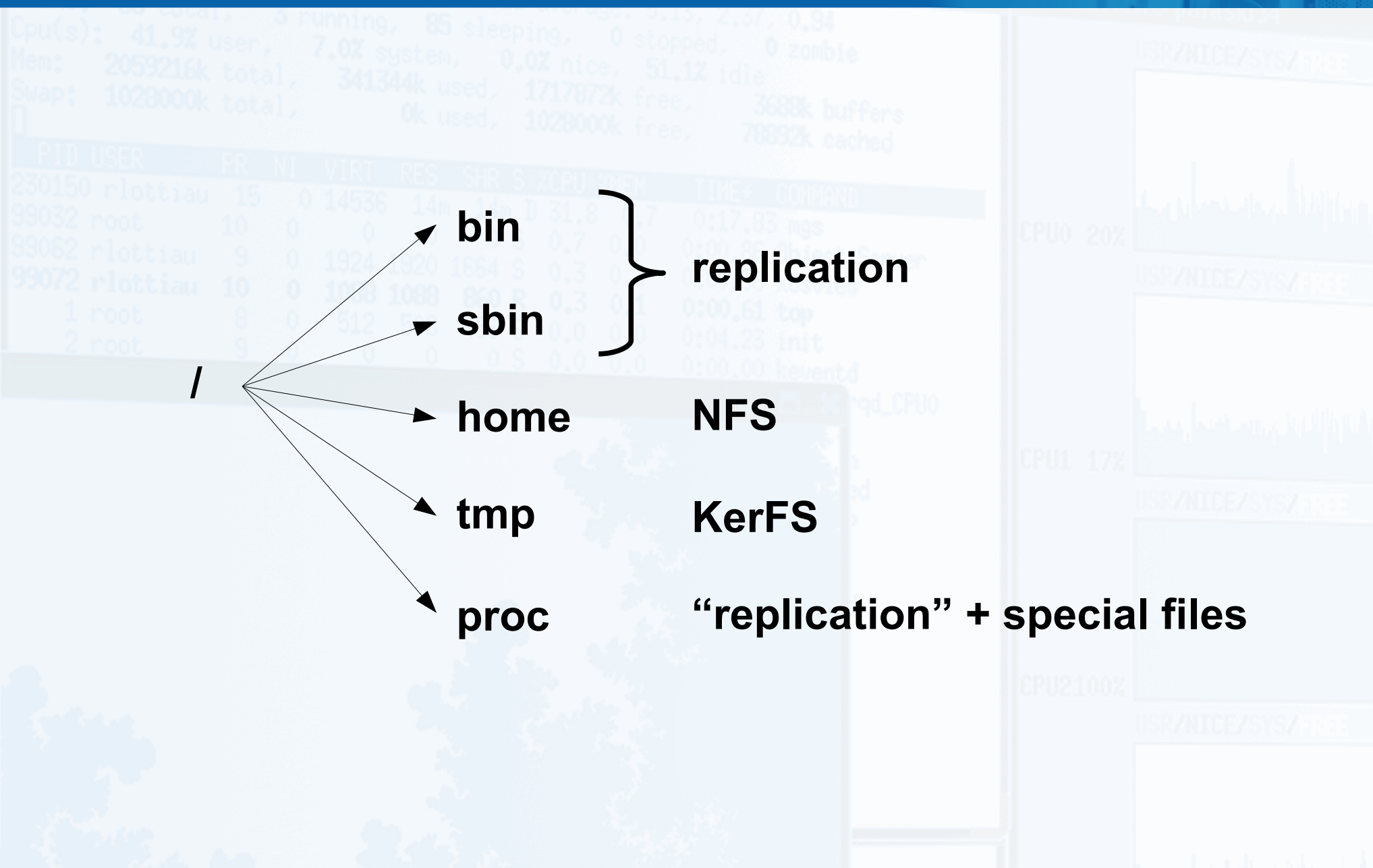
- In Kerrighed processes can migrate
  - After a migration, processes must see same files and same file content
  - All processes must see same files and same content
- We (at least Kerlabs) need strong and stable solution
  - Customers would not accept unstable or experimental FS
  - Many customer will ask for a specific file system

/ →
- bin ⎫
- sbin ⎭ replication
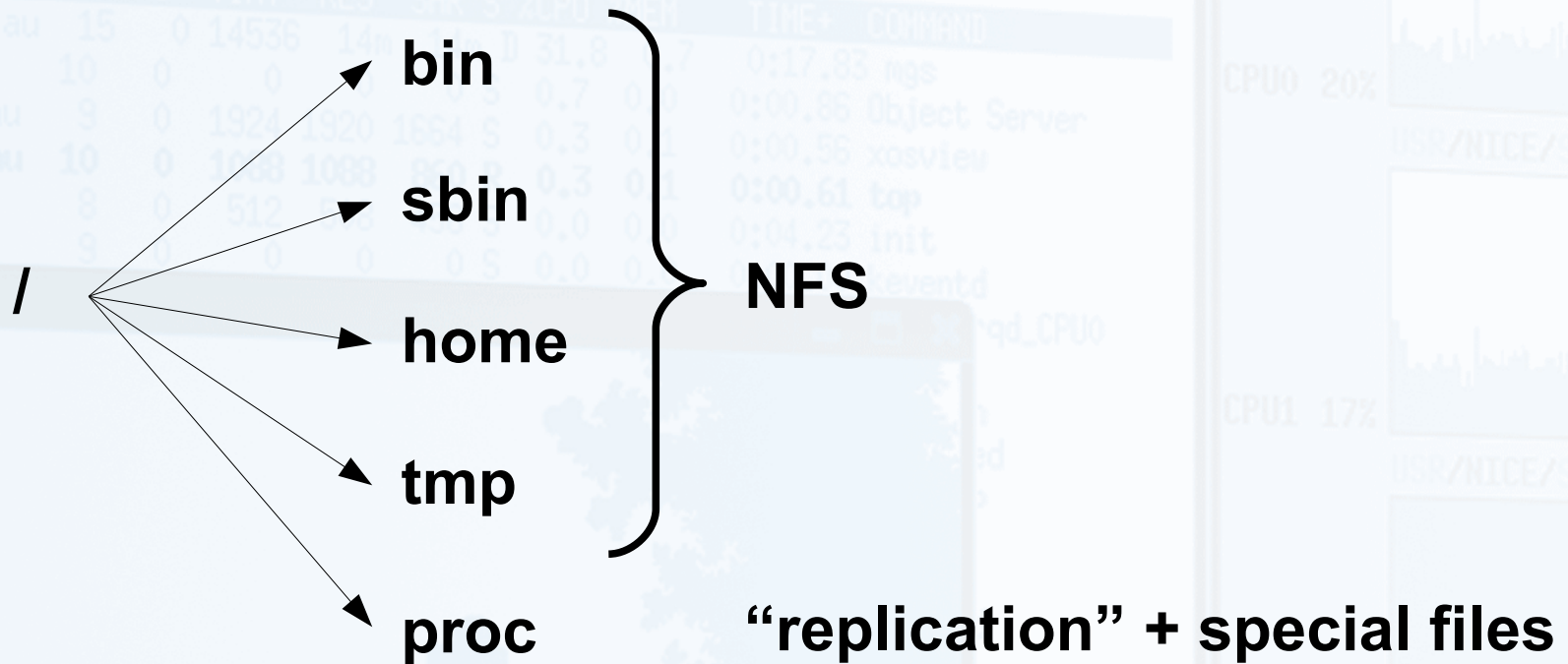- home → NFS
- tmp → KerFS
- proc → "replication" + special files

- "Simple" but that simply works !

```
        bin

        sbin
                          NFS
/       home

        tmp

        proc        "replication" + special files
```

- Using only NFS root
  - We have no choice to offer
  - Performance could be an issue
  - Installing an NFS-root server is not obvious
- Would be nice to
  - At least offer other choices
  - Offer a really integrated FS for Kerrighed

**KERLABS**

# Offering other choices

- We have plenty of file system to test
  - Free systems
    - GFS
    - OCFS
    - NFS V4
      - NFS V4.1 for parallel NFS
    - Lustre (probably too complex)
    - CIFS !
    - Etc…
  - Proprietary systems
    - SFS (SeaNodes FS + Hexanode)
    - Isilon solutions
    - Etc…
- Quite a lot of work to check all these FS

# Issue #1 : Existing DFS are not Kerrighed aware

- **Cluster wide coherence of file data**
  - What happens if 2 processes on 2 different nodes access the same file at the same time ?
    - NFS : data is not coherent during a short period of time
    - GFS : data is coherent, but the coherence is expensive

- **File access performance**
  - NFS : poor
  - GFS : good if good file access locality

- **Unix sockets does not work cluster wide !**
  - This is not a file system problem
    - Solved with the migrable streams

Kerrighed

KERLABS

# Offering an integrated file system

- That was the plan for KerFS

- Now we have the kDFS file system under implementation
  - Storage of data on different cluster nodes
  - Implementation "Kerrighed aware"
    - Use of KDDM
    - No bad effect with process migration
    - Possible to implement intelligent policies for data / process placement

KERLABS

- Designing a production quality new FS from scratch is a very hard task
  - A whole company could be needed to achieve this
- The big question is
  - Would it be possible to reach a production quality FS ?
  - Yet another research FS is fun !
    - But it's just another fun research FS...
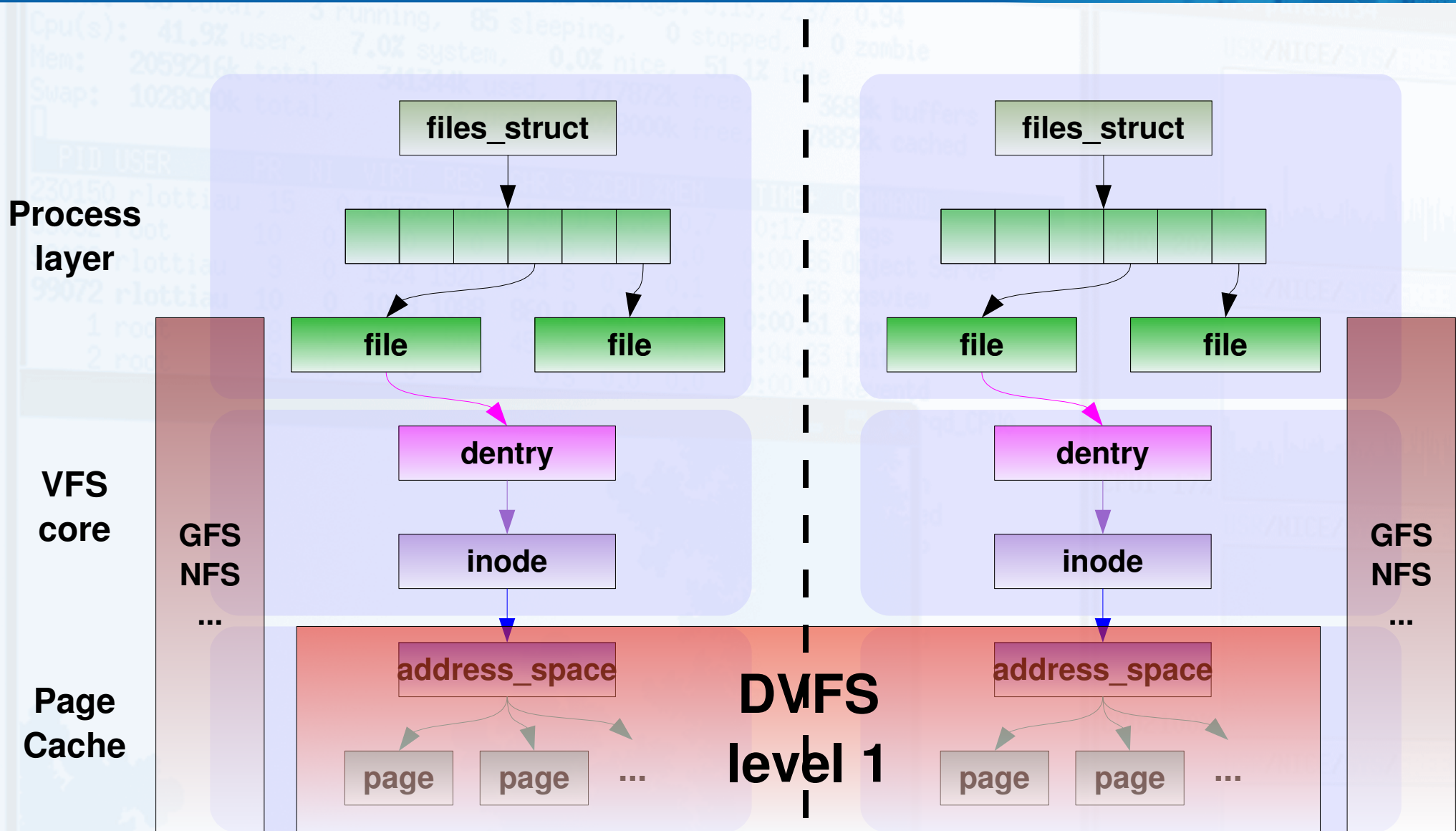- If we consider node failures, the problem is even more complicated

# A third approach: DVFS

- The main problem of existing DFS : data coherence

- Why there is no coherence problem without a DFS ?
  - No distributed copies of data : one unique file cache

- Solution : 1 unique cluster wide file cache !
  - Kerrighed DVFS level 1

- Data has no more to be written back to the server or the shared disk to be accessed remotely

- Data can be accessed directly from one node memory to another node memory

**Process layer**
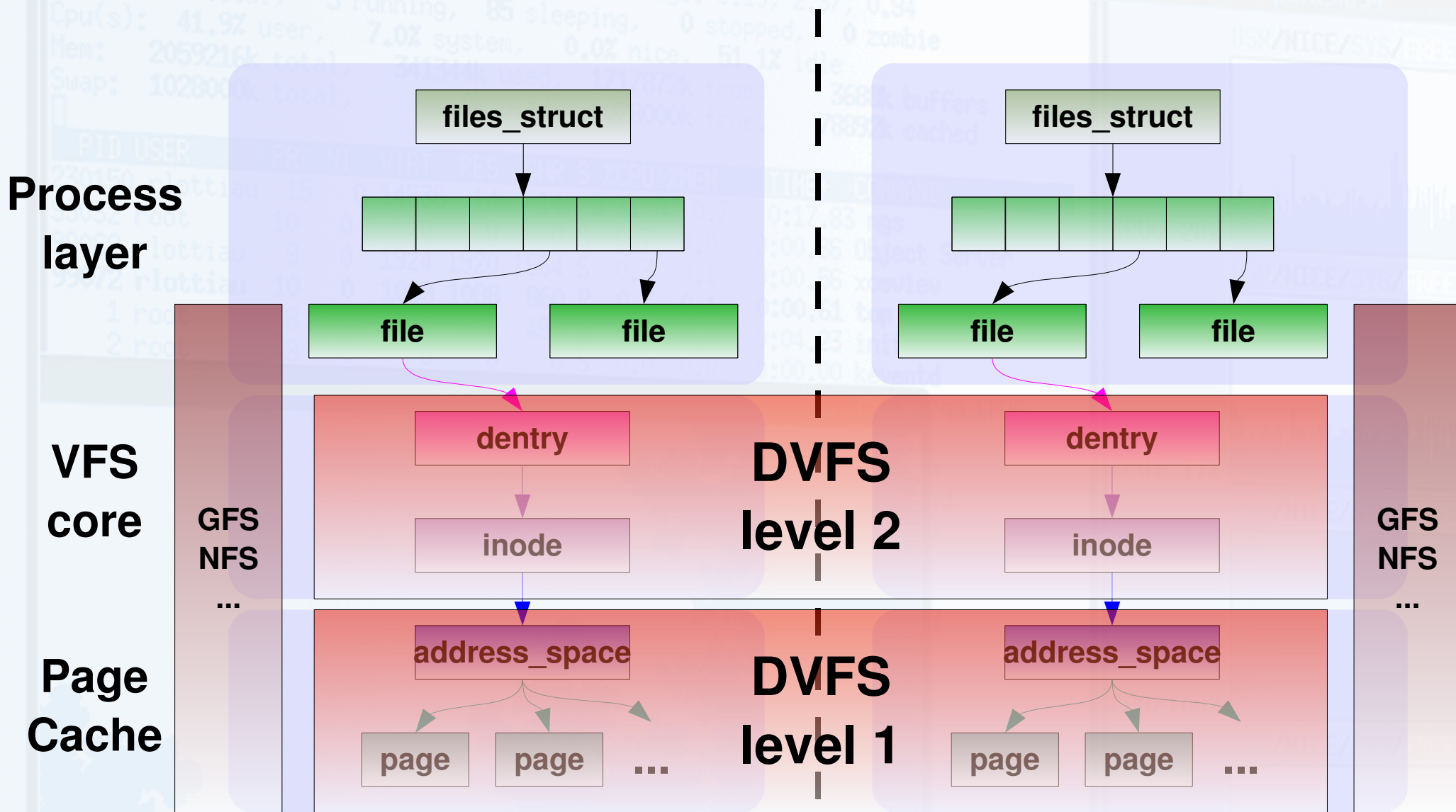
files_struct

file | file

**VFS core**

GFS NFS ...

dentry

inode

**Page Cache**

GFS NFS ...

address_space

page | page | ...

**DVFS level 1**

files_struct

file | file

dentry

inode

address_space

page | page | ...

Kerrighed

KERLABS

**Process layer**

files_struct

file — file

**VFS core**

GFS NFS ...

**DVFS level 2**

dentry

inode

**Page Cache**

GFS NFS ...

**DVFS level 1**

address_space

page — page — ...

Kerrighed

KERLABS

# Linux VFS – Kerrighed DVFS level 3

**Process layer**

**VFS core**

**Page Cache**

files_struct

files_struct

DVFS level 3

file · file

file · file

dentry

dentry

DVFS level 2

inode

inode

GFS NFS ...

GFS NFS ...

address_space

address_space

DVFS level 1

page · page · ...

page · page · ...

- Kerrighed aware architecture

- Can be plugged to "any" existing FS

- New DFS can be implemented within this framework

- Simpler to implement than a all new brand DFS

- This is not a new FS

  - This is more acceptable for customers

- Nothing done yet

- No manpower to do that

- Simpler than a DFS but quite complex anymay