



Owncloud scalability and a Nextcloud design for 10.000-20.000 users.



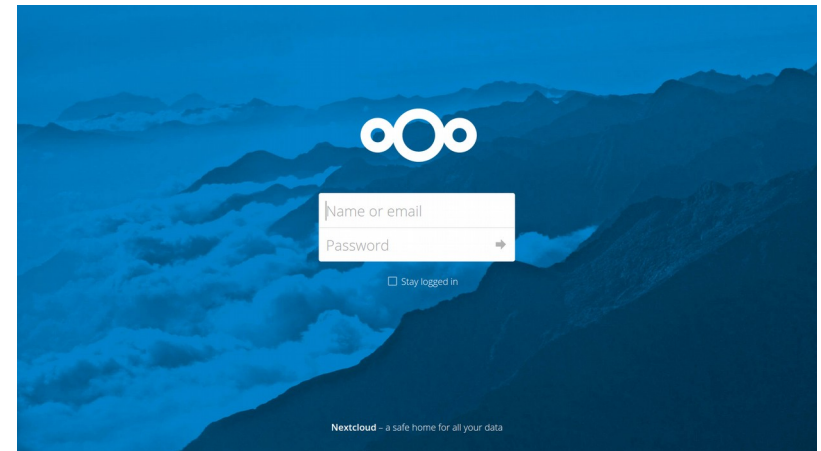
# Introduction

- Dennis Pennings
- 360 ICT (.nl)



## The goals

Design a 20.000 user  
NC implementation.



# Documentation

(docs.nextcloud.com)

## Recommended System Requirements

4 to 20 application/Web servers.

A cluster of two or more database servers.

Storage is an NFS server, or an object store that is S3 compatible.

Cloud federation for a distributed setup over several data centers.

Authentication via an existing LDAP or Active Directory server, or SAML.

**100.000users / 1Pb**

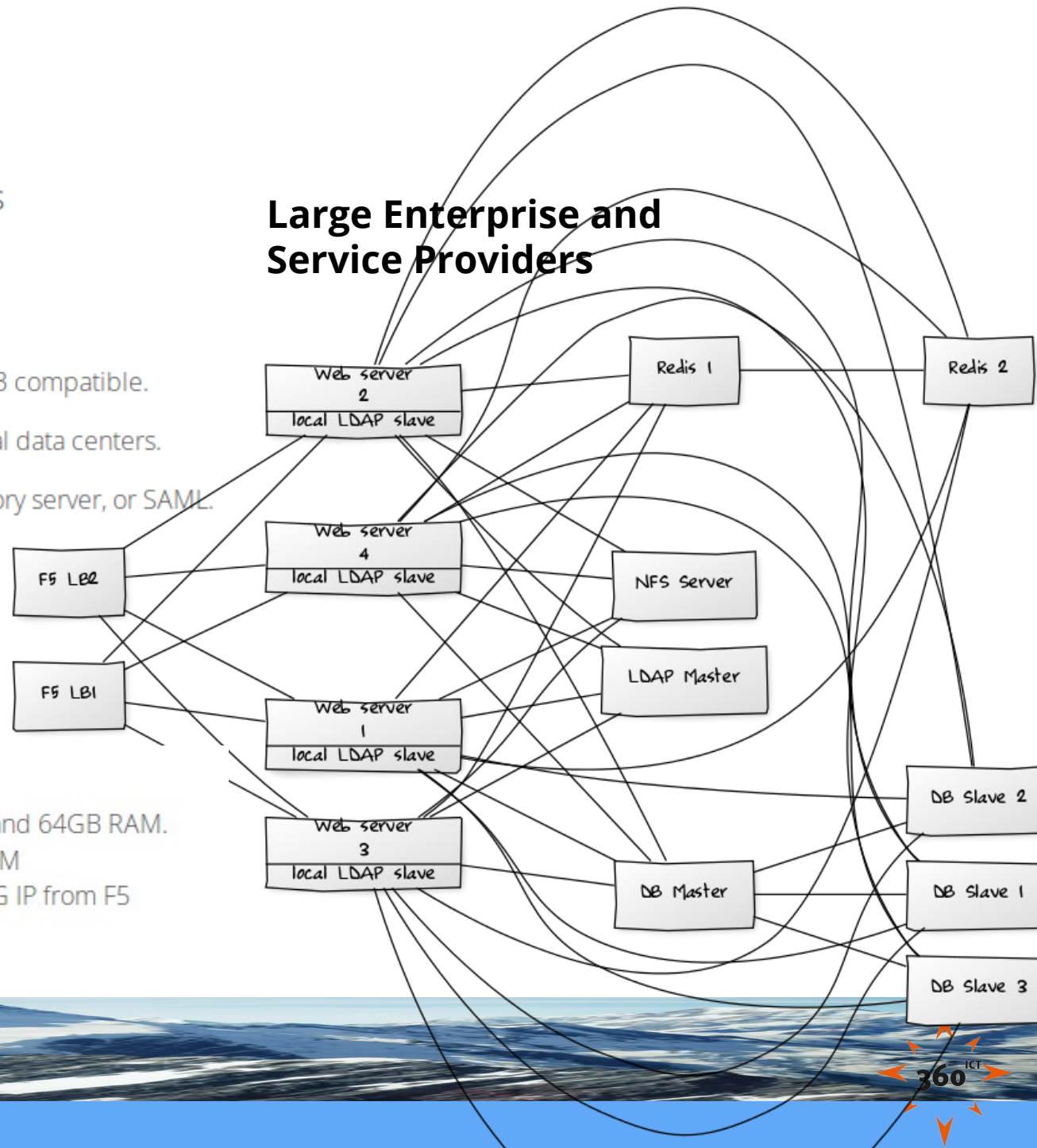
**Database: MySQL/MariaDB Galera  
Cluster with 4x master - master  
replication.**

### Multisite?

#### • Components

- 4 to 20 application servers with 4 sockets and 64GB RAM.
- 4 DB servers with 4 sockets and 128GB RAM
- 2 Hardware load balancer, for example BIG IP from F5
- NFS storage server as needed.

## Large Enterprise and Service Providers



# Documentation

([docs.nextcloud.com](https://docs.nextcloud.com))

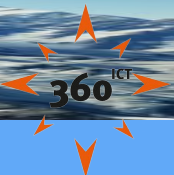
**Large Enterprise and  
Service Providers**

- Additional notes
  - Use LDAP slaves on webserver.
  - Use SSL offloading on load balancers.
  - Redis for session management storage.
  - Redis for in-memory caching.
  - Redis: provides persistence, nice graphical inspection tools available, supports Nextcloud high-level file locking.
  - Memcached if shibboleth is used.



# (Somewhat) Large Deployments

- Info collected from CS3 2016  
[cs3.ethx.ch](http://cs3.ethx.ch)
- Presented on Nextcloud conference with a Q&A from Frank.
  - [360ict.nl/blog](http://360ict.nl/blog)
  - [youtube.com/nextcloud](https://youtube.com/nextcloud)
- All data in a sheet on [360ict.nl/blog](http://360ict.nl/blog) and I will update it as I receive new information.



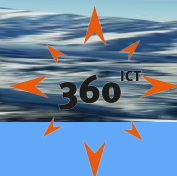
# (Somewhat) Large Deployments

| Name                             | <u>CERNbox</u>    | <u>SURFdrive</u>       | <u>u:cloud</u>          | <u>MyCore</u>         | <u>SWITCHdrive</u>                 |
|----------------------------------|-------------------|------------------------|-------------------------|-----------------------|------------------------------------|
| Date of inquiry                  | Jan 2016          | Jan 2016               | Jan 2017                | Jan 2016              | Dec 2016                           |
| <b>Users</b>                     |                   |                        |                         |                       |                                    |
| total users                      | 5k                | 13k                    | 10k                     | 4k-15k (jan-dec)      | 23k                                |
| max number of users              |                   |                        |                         |                       |                                    |
| concurrent per day               |                   | 4000                   | 3600 seen, 2000 average |                       | 10k per week                       |
| number of files                  | 55m (500m on EOS) | 75m                    | 5M                      |                       | 60m                                |
| <10Mb                            |                   | 95%                    | 98%                     |                       |                                    |
| <1Mb                             |                   | 90%                    | 85%                     |                       |                                    |
| data over multiple users divided |                   |                        | 2000 users have 10Tb    |                       |                                    |
| quota per user                   | 1Tb               | 8Gb                    | 50Gb                    | 20Gb                  | 50Gb                               |
| Linux:mac:win %                  |                   |                        |                         | 20-20-40              |                                    |
| costs                            |                   | 18eu per user per year |                         |                       |                                    |
| <b>Database</b>                  |                   |                        |                         |                       |                                    |
| type                             | none              | MySQL                  | MySQL                   | MariaDB               | PostGress → mariadb                |
| setup                            | none              | Galera                 |                         | Galera-clustercontrol | Galera                             |
| number of nodes                  | none              | 4 physical nodes       | 1                       | 3 VMs                 | 4 VMs                              |
| cpu/mem                          | none              | 32 cores, 256Gb        | 2 cores, 12 Gb          | 8 cores, 16Gb         | 46 vCPU, 250Gb                     |
| storage                          | none              |                        | 200GB                   | local SSD storage     | local SSD                          |
| size of database                 |                   |                        |                         |                       |                                    |
| <b>Network</b>                   |                   |                        |                         |                       |                                    |
| LB engine webserver              |                   |                        |                         |                       |                                    |
| LB engine webserver              |                   |                        | F5                      | 2 VMs, 4 cores, 2Gb   | Haproxy, 4 vCPU / 4GB              |
| LB engine database               |                   | Haproxy → maxscale     |                         | Haproxy               |                                    |
| Bandwidth Inet/DB/storage        |                   |                        |                         |                       |                                    |
| Inet bandwidth used              |                   |                        |                         |                       | extremely variable, peaks at 8Gb/s |



# (Somewhat) Large Deployments

| Name                             | Cloudstor                                     | Sciebo  | UniPiBox | Polybox                        |
|----------------------------------|---|---|----------|--------------------------------|
| Date of inquiry                  | Jan 2017                                      | Dec 2016  | Jan 2016 | Jan 2017                       |
| <b>Users</b>                     |   |   |          |                                |
| total users                      | 30k   | 25k   |          | 21000                          |
| max number of users              |   | 500k  |          |                                |
| concurrent per day               | 3k per day                                    | 7000  | 75       | 5000-6000                      |
| number of files                  | 43m   | 20m   | 350k     | 35m                            |
| <10Mb                            | 96%   |   |          | 99%                            |
| <1Mb                             | 82%   |   |          |                                |
| data over multiple users divided |   |   |          |                                |
| quota per user                   | 100gb,<br>1Tb per project (change on request) | 30Gb users, 500Gb for employees,<br>30Gb-2Tb for projects | 10Gb/3Gb |                                |
| Linux:mac:win %                  | 53-7-40                                       |   |          |                                |
| costs                            |   |   |          |                                |
| <b>Database</b>                  |   |   |          |                                |
| type                             | Mariadb                                       | MariaDB   |          | Percona XtraDB Cluster (MySQL) |
| setup                            | Galera  | Galera  |          | Galera                         |
| number of nodes                  | 12  | 4 physical nodes  |          | 3 VM's active, 4th on standby  |
| cpu/mem                          | 32 cores, 256Gb                               | 20 cores, 256Gb   |          | 16 core, 40Gb                  |
| storage                          | local SSD                                     | 800Gb SSD with raid10                                     |          | 100Gb SSD per node             |
| size of database                 |   |   |          | 35Gb                           |
| <b>Network</b>                   |   |   |          |                                |
| LB engine webservers             |   | LVS with keepalive  |          | LVS with keepalive             |
| LB engine webservers             | HA proxy                                      | 8 core Linux machines per site                            |          |                                |
| LB engine database               | maxscale                                      | Maxscale  |          | Maxscale                       |
| Bandwidth Inet/DB/storage        | 120Gb L3VPN                                   | 10Gb/10Gb/56Gb  |          | 10Gb/10Gb/10Gb                 |
| Inet bandwidth used              |   | 80MB/sec  |          | 6.2MB/s (3.5Tb per week)       |



# (Somewhat) Large Deployments

| Name                    | Date of inquiry | CERNbox<br>Jan 2016 | SURFdrive<br>Jan 2016                                 | u:cloud<br>Jan 2017               | MyCore<br>Jan 2016               | SWITCHdrive<br>Dec 2016  |
|-------------------------|-----------------|---------------------|---|-----------------------------------|----------------------------------|--|
| <b>Storage</b>          |                 |                     |   |                                   |                                  |  |
| <u>number of nodes</u>  |                 | 1300                | 9   |                                   | 2 arrays, each with 6 nodes      |  |
| <u>number of disks</u>  |                 | 40k disks           |   |                                   |                                  |  |
| <u>cpu/mem</u>          |                 |                     |   |                                   | 8 cores, 12Gb, 800Gb SSD         |  |
| <u>brand/type</u>       |                 |                     |   |                                   | Dell PowerEdge R620/630/MD3420   | No-name  |
| <u>software</u>         |                 | EOS                 | GlusterFS (distr-rep)                                 | Object storage, scalability, Fuse | Scality                          | nfs on top of Ceph   |
| <u>Total space</u>      |                 | 1,3PB (64PB on EOS) | 100Tb   | 1.8Pb                             |                                  | 1.6Pb  |
| <u>space in use</u>     |                 | 104TB               |   | 10Tb (OC)                         |                                  | 58Tb   |
| <u>Issues</u>           |                 |                     | add server takes 2 months to rebalance, Backup issues |                                   |                                  | after 4k users, Db to local SSD, ceph volumes from 100Tb lazy to 2Tb |
| <u>Futureplans</u>      |                 |                     | GPFS, EOS, dCache                                     |                                   |                                  | object storage, Qowncloud sharding with Federation between instances |
| <b>Webserver</b>        |                 |                     |   |                                   |                                  |  |
| <u>number of nodes</u>  |                 |                     | 12  | 1                                 | 6 VMs                            | 10 sync / 4 web  |
| <u>cpu/mem</u>          |                 |                     |   | 4 cores, 20Gb                     | 8 cores, 24Gb                    | Sync: 16GB, 16vcpu / Web: 4Gb, 4vcpu                                 |
| <u>webengine</u>        |                 |                     | Apache  | Apache                            |                                  | Apache   |
| <u>php version</u>      |                 |                     |   | 5.6                               |                                  | 5  |
| <u>OC version</u>       |                 | 8                   | 8.2   | 8.23                              | 8.2                              | 8.1  |
| <b>other info</b>       |                 |                     |   |                                   |                                  |  |
| <u>Issues</u>           |                 |                     |   |                                   | Version app causes too much load |  |
| <u>OS</u>               |                 |                     |   | CentOS                            |                                  |  |
| <u>identity</u>         |                 |                     |   | Shibboleth + LDAP                 | shibboleth                       | ldap   |
| <u>Uses docker</u>      |                 | no                  | no  | no                                | no                               | yes  |
| <u>management tools</u> |                 |                     |   |                                   |                                  | ansible  |
| <u>Other</u>            |                 |                     |   |                                   |                                  |  |



# (Somewhat) Large Deployments

| Name                    | Cloudstor   | Sciebo   | UniPiBox  | Polybox        |
|-------------------------|---|--|---|----------------|
| Date of inquiry         | Jan 2017  | Dec 2016   | Jan 2016  | Jan 2017       |
| <b>Storage</b>          |   |  |   |                |
| <u>number of nodes</u>  | 12 nodes  | 10 nodes   | 2X RADOS-GW, 8 CORES, 16GB, 2X KEYSTONE, 2 CORES 1GB          |                |
| <u>number of disks</u>  | 504   |  |   |                |
| <u>cpu/mem</u>          | Xeon E5, 128GB  | 20 cores, 256Gb  |   |                |
| <u>brand/type</u>       | supermicro  | IBM GSS appliance  |   | IBM SoNAS      |
| <u>software</u>         | EOS   | GPFS   | Ceph, radosgw+keystone for swift                              | NFS            |
| <u>Total space</u>      | 1.9PB   | 5PB  | 2PB   |                |
| <u>space in use</u>     |   | 800Tb (including snapshots)  |   | 62TB           |
| <u>Issues</u>           |   | none   |   | none           |
| <u>Future plans</u>     | Cernbox   |  | 30 SITES, DNS-ha AS lb, Ceph, , xtradb cluster, swift storage |                |
| <b>Webserver</b>        |   |  |   |                |
| <u>number of nodes</u>  |   | 16 physical nodes  | 2 VMs   | 8 VMs          |
| <u>cpu/mem</u>          | 48  | 16 cores, 128Gb  | 4 cores, 16GB   | 16 core, 16Gb  |
| <u>webengine</u>        | Apache  | Apache   |   | Apache         |
| <u>php version</u>      | 5.6   | 7  |   | 5.4            |
| <u>OC version</u>       | 8.2   | 9.06   |   | 8              |
| <b>other info</b>       |   |  |   |                |
| <u>Issues</u>           | Database design is rubbish, inefficient query and query caching, redundancy of metadata, poor group sharing design, flawed assumptions and key selections, not built to scale |  |   |                |
| <u>OS</u>               | RHEL7.3   | Redhat 6/7   |   | RHEL7          |
| <u>identity</u>         | SAML  | Shibboleth for registration, LDAP to log in to OC  |   |                |
| <u>Uses docker</u>      | yes   | no   | no  | no             |
| <u>management tools</u> | ansible, rancher, jenkins   | Ansible  | puppet  | Zabbix, Splunk |
| <u>Other</u>            | Cluster spans a 65ms network end to end. Looking at moving to Cernbox to avoid the database issue.  | Users distributed to three sites, backup, but no load balancing between sites. One OC instance per institution |   |                |

# (Somewhat) Large Deployments

- A summary (mostly cs3-2016 info)
  - *8 different implementations. 5k-30k unique users. 2k-10k concurrent users.*
  - *Mostly 4 node MySQL-MariaDB Galera with maxscale. 2-46 cores, 12-256Gb mem. One node which handles writes (with failover).*
  - *Webservers, Apache, php5, OC8/OC9, 1-48 cores, 16-128Gb mem.*
  - *100Tb-800Tb in use, 1.3-5Pb allocated.*
  - *A few use docker.*
  - *Scaling problems/limit with DB.*

Not much differences here.





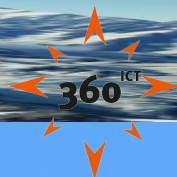
# (Somewhat) Large Deployments

- Storage looks very different
  - *GlusterFS*
  - *Ceph*
  - *Ceph-NFS*
  - *NFS*
  - *Scality-Fuse*
  - *EOS*



# Initial Concept Design

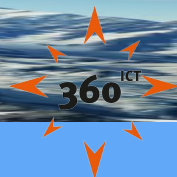
|                  |                     | NextCloud   | 360 ICT concept design                |
|------------------|---------------------|---|---------------------------------------|
| <b>Users</b>     | total users         | 5k  | 10k                                   |
|                  | max number of users | 100k  | 20k                                   |
|                  | concurrent per day  |   |                                       |
|                  | number of files     |   |                                       |
|                  | quota per user      |   | 5Gb                                   |
| <b>Webserver</b> | number of nodes     | 4 to 12   | 4-n VM's                              |
|                  | cpu/mem             | 4 sockets / 64Gb                                  | 1-4 core, 10Gb-20Gb                   |
|                  | webengine           | Apache  | Apache                                |
|                  | php version         | 5.5+  | 7                                     |
|                  | OC version          | Enterprise Edition                                | NC11 (because of the S3 improvements) |
| <b>Database</b>  | type                | MySQL/MariaDB<br>(Oracle/PostgreSQL is supported) | MariaDB                               |
|                  | setup               | Galera  | Galera                                |
|                  | number of nodes     | 4   | 4 VM's                                |
|                  | cpu/mem             |   | 4-24 cores / 24-256 Gb                |
|                  | storage             | SSD   | SSD                                   |





# Initial Concept Design

|                   | NextCloud        | 360 ICT concept design                  |
|-------------------|------------------|---|
| <b>Network</b>    | LB engine        | Kemp VLM or maxscale                    |
|                   | LB hardware      | Kemp VM / appliance or VM with maxscale |
| <b>Storage</b>    | number of sites  | 1 to 2                                  |
|                   | number of nodes  | 4-n nodes                               |
|                   | cpu/mem          | 4-8 cores,                              |
|                   | brand/type       | Dell PowerEdge                          |
|                   | software         | Swift (openstack or ceph)               |
| <b>Other info</b> | Total space      | 50Tb to 100Tb                           |
|                   | OS               | Ubuntu 16.04                            |
|                   | identity         | MS AD or Shibboleth                     |
|                   | Uses docker      | yes                                     |
|                   | management tools | openstack, swarm                        |



# Q&A for our initial design

- Some highlights from the Nextcloud Q&A session
  - *Multiple datacenters will be hard because the dependency of DB and storage which need to be in sync.*
  - *Docker supported by Nextcloud (OwnCloud does not)*
  - *No preference for database, mysql/mariadb are mentioned as a starting point, but others are supported (like Oracle). That's true for Loadbalancers too (haproxy/maxscale/hardware).*
  - *Nextcloud has no preference on OS, any version with (possible) enterprise support is supported*





# Q&A for our initial design

- Some highlights from the Nextcloud Q&A session
  - *Object storage is supported, multibucket storage available in NC11/OC9.*
  - *Use php v7 if possible if your distro supports it.*
  - *There's definitely a scale limit to the database. NC11 claims better database scalability (bold statements like 80% less queries)*
  - *Storage is a concern, but in all fairness this is a client specific issue. OC/NC support lots of solutions.*
  - *The full session from the Q&A session (1 hour) can be found on [360ict.nl/blog](http://360ict.nl/blog) and [youtube.com/nextcloud](https://youtube.com/nextcloud).*



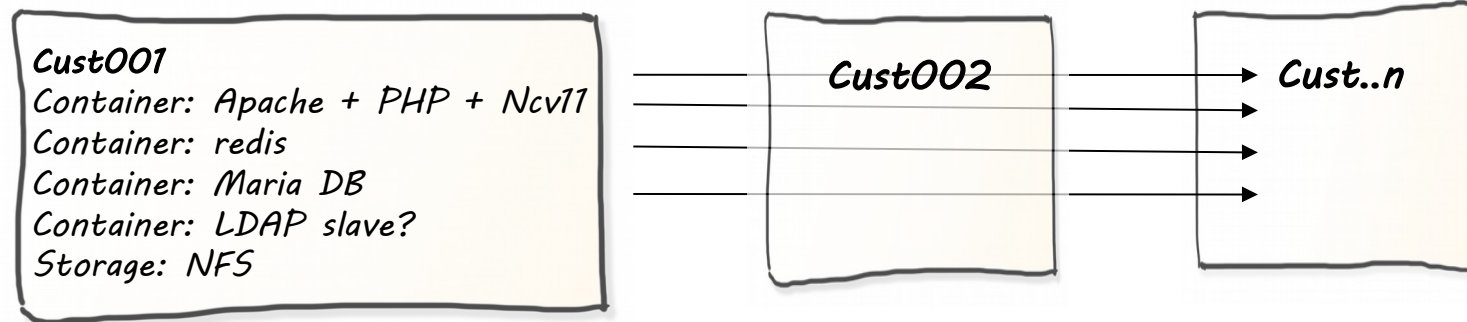
# Our current design

- Lots of small instances instead of 1 (or more) monolithic approach.
- Based on Docker containers.
- Kubernetes as container orchestration
- Distributed storage abstracted in the hypervisor. NFS storage offered through VMs on the hypervisor.



# Our current design

## Next Cloud Pod

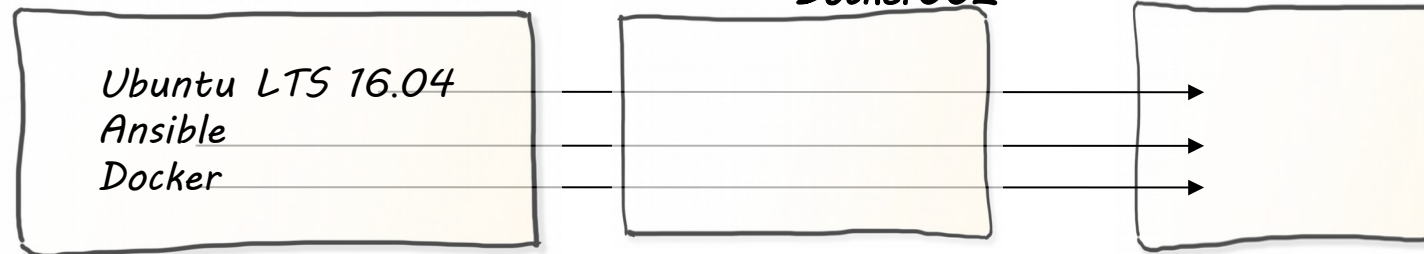


## Host VM

## Docker001

## Docker002

## Docker..n

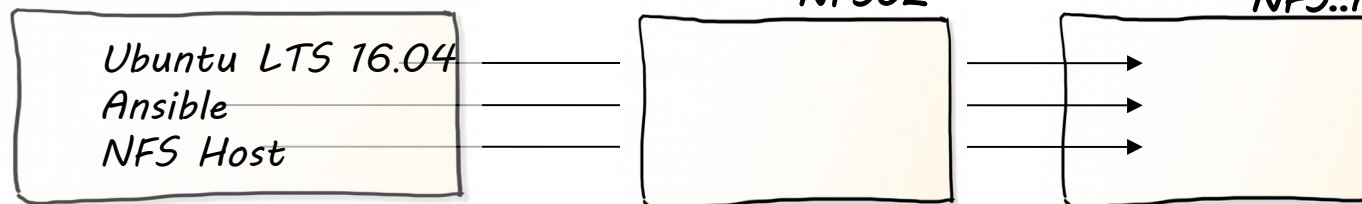


## Storage VM

## NFS01

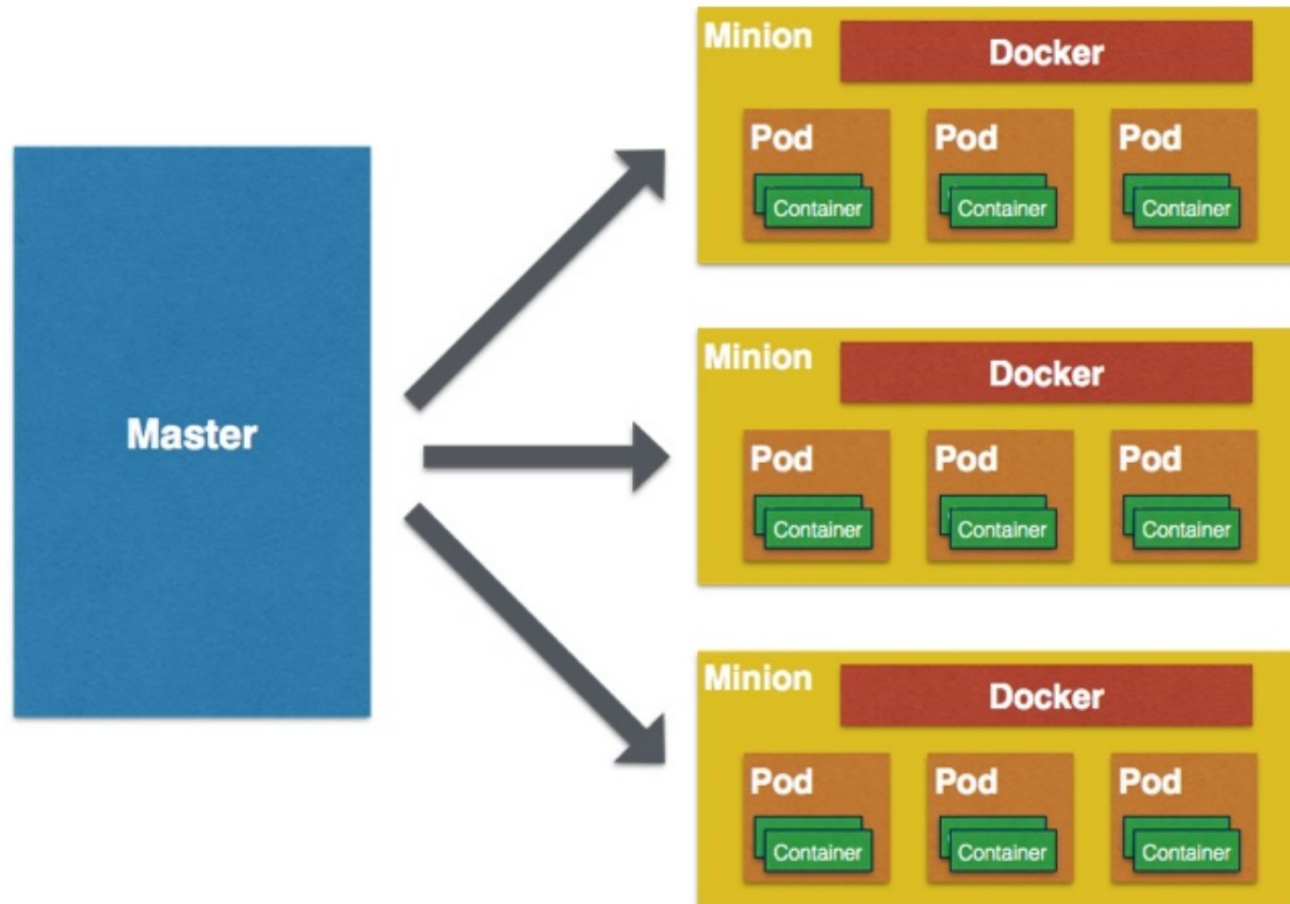
## NFS02

## NFS..n

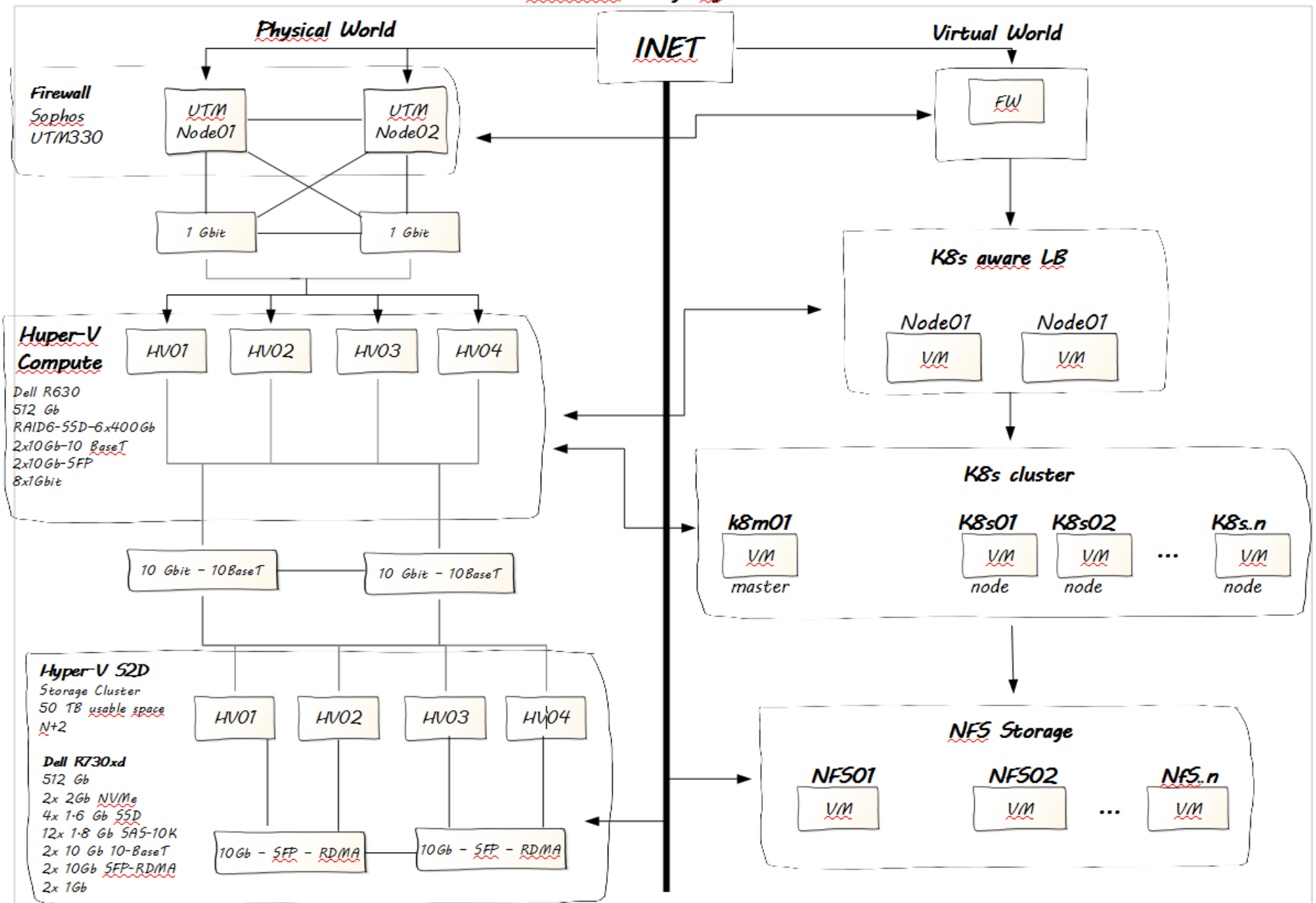




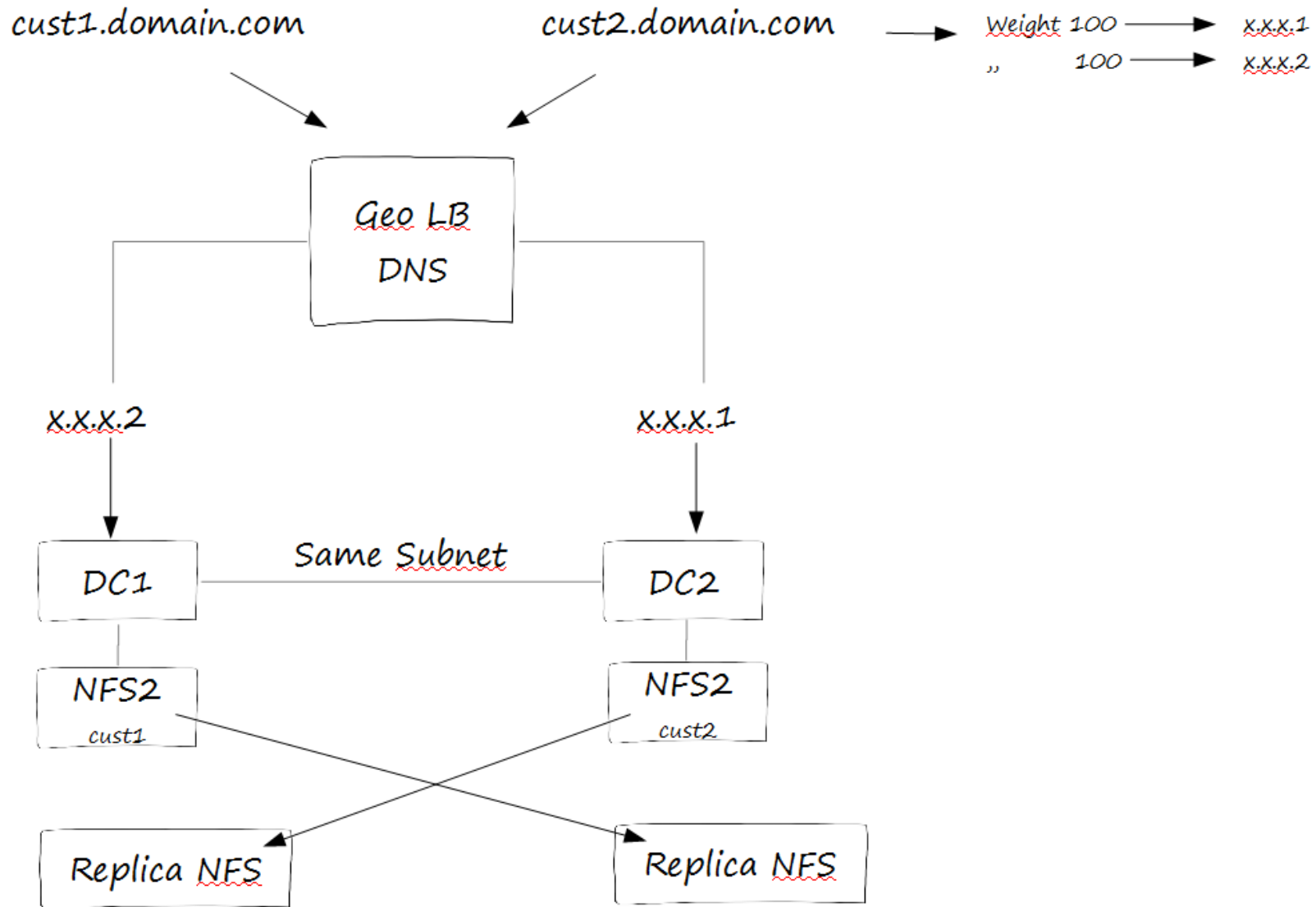
# Kubernetes



<http://www.slideshare.net/roland.huss/fabric8-and-docker-kubernetes-openshift>



# Multiple Datacenters





# Our concept design

- Advantages

- *Simple setup of NC. U:cloud can scale to thousands of users with 2 servers. But our typical users per instance will be 1-20 users with a few instances of 100's of users.*
- *Customization per customer is possible (but maybe not wanted?).*
- *Updating can be done per instance/customer. So rollout is gradually and not all or nothing.*
- *No need for large (Galera) clusters or web/db load balancing.*



# Our concept design

- Disadvantages

- *NFS storage is unavailable during reboots. NC will be paused if NFS is unavailable and will recover. Not ideal, but good enough for first steps.*
- *Performance of Virtualized NFS for storage? NFS will probably be an in-between step.*
- *Will it scale? We are planning loadtests in Google Cloud Engine and our own hardware. We are curious as how far we can push the design, the max number of users in a pod and also tests on different kind of storages (and how much it would cost to do a 100k user test in GCE).*



# Our concept design

- Other stuff
  - *Use federation for larger instances that don't fit in an pod?*
  - *The backend in this design is based on MS distributed storage (S2D), but this can be any storage: distributed (GlusterFS, Ceph, Netapp, etc), (small) SANs or fileclusters.*
  - *The storage can be split up in this design, which allows for more flexible solutions.*
  - *Maybe even running GlusterFS as containers within kubernetes?*
  - *Nextcloud GS?*







# Questions?



If you have a large implementation of OC/NC and want to share your experiences, email me at **[dennis@360ict.nl](mailto:dennis@360ict.nl)**.  
The sheet will be shared on [360ict.nl/blog](http://360ict.nl/blog)