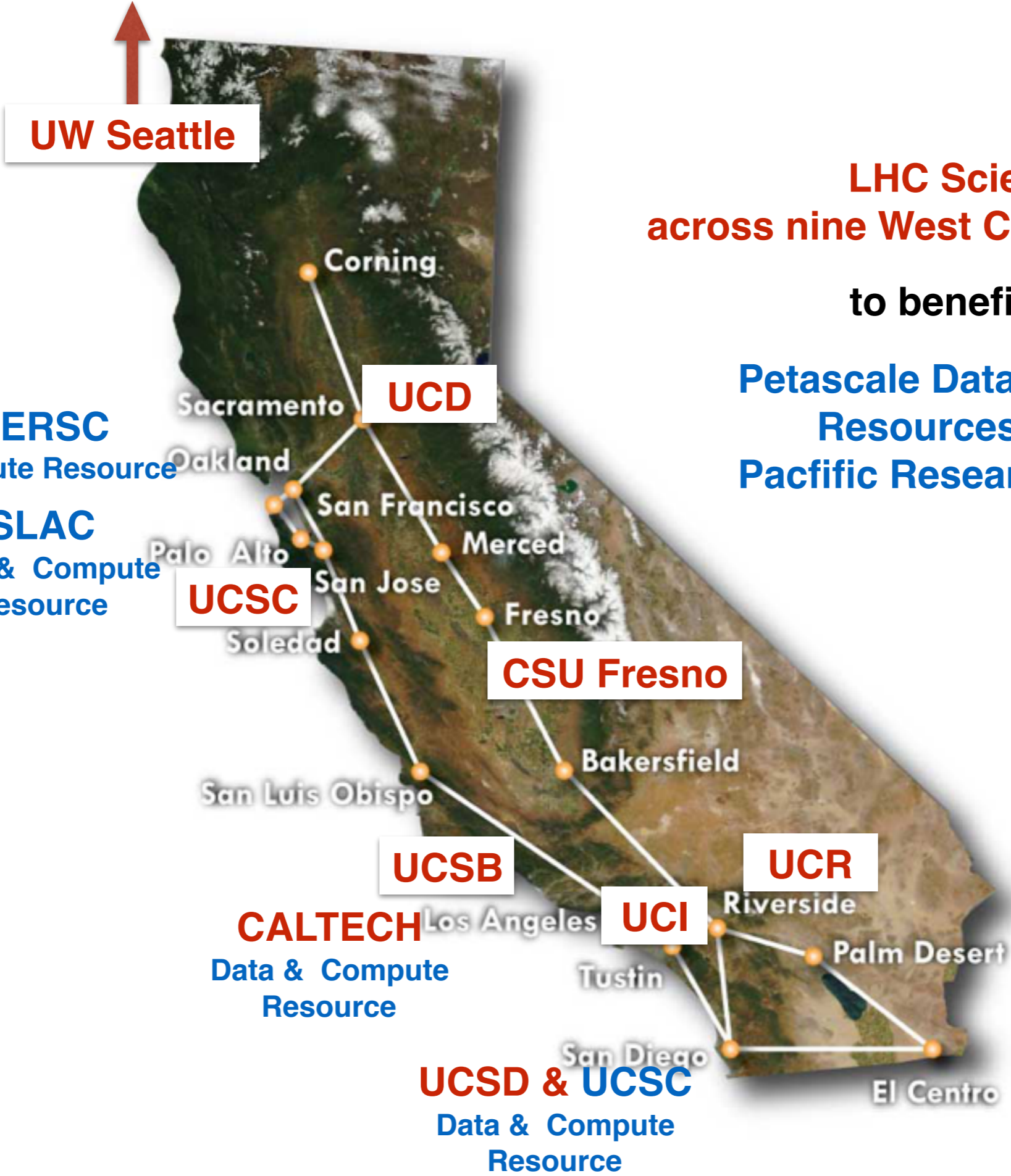


Flying HTCondor at 100gbps Over the Golden State

Jeff Dost (UCSD)

What is PRRP?

- Pacific Research Platform:
 - 100 gbit network extending from Southern California to Washington
 - Interconnects Science DMZ between institutions



**LHC Scientists
across nine West Coast Universities**

to benefit from

**Petascale Data & Compute
Resources across
Pacific Research Platform**

**NERSC
Compute Resource**

**SLAC
Data & Compute
Resource**

**CALTECH
Data & Compute
Resource**

**UCSD & UCSC
Data & Compute
Resource**

What is LHC @ UC?

- Pilot project that unifies University of California Tier 3 campus compute infrastructures on top of PRP network
- Provides each UC the ability to:
 - Utilize external compute resources
 - Provided by GlideinWMS HTCondor pool
 - Access data from anywhere
 - Provided by XRootD
- Integrates seamlessly with local compute and storage resources

What is LHC @ UC?

- Participating UCs:
 - ATLAS
 - CMS
- Resources currently provided:
 - Each UC
 - 50k core Comet cluster at SDSC
- Eventually:
 - Any other non-UC participating PRP site
 - Any OSG site beyond PRP
 - Other NSF XSEDE and DOE super computing centers
 - Cloud resources



Data Access

- Built on top of XRootD
- Jobs don't need to run where the data is located
- Local UC and external compute resources both cache remote data accesses
- Arbitrary data from local UC can be exported and made available to all compute resources

Hardware shipped to UCs

(aka the “brick”)



Hardware:

- 40 cores
- 12 x 4TB data disks
- 128 GB ram
- 2 x 10 gbit network interface

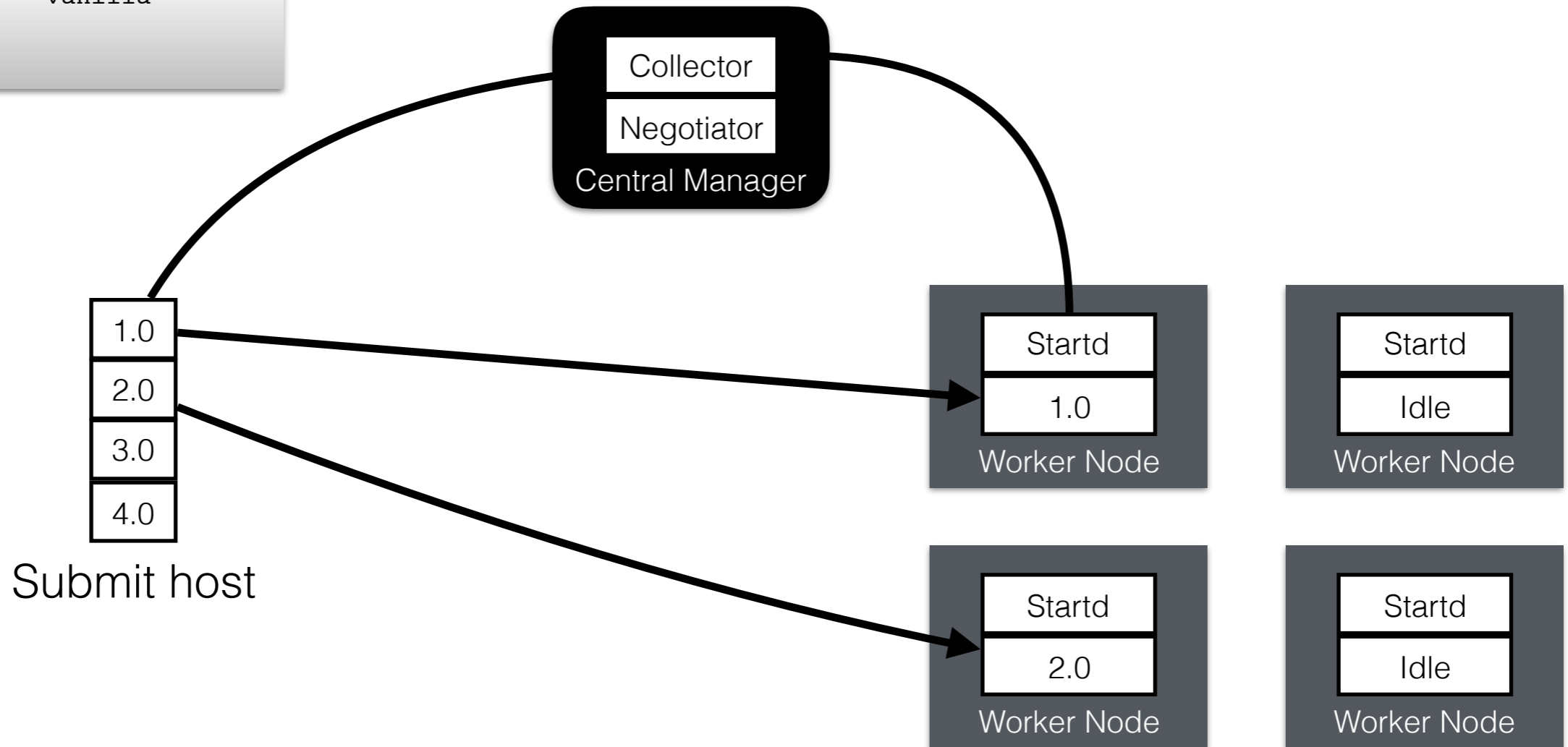
Software:

- Full HTCondor pool
- XRootD server, redirector, and proxy cache
- cvmfs w/ optional Squid

The brick is effectively a site in a box

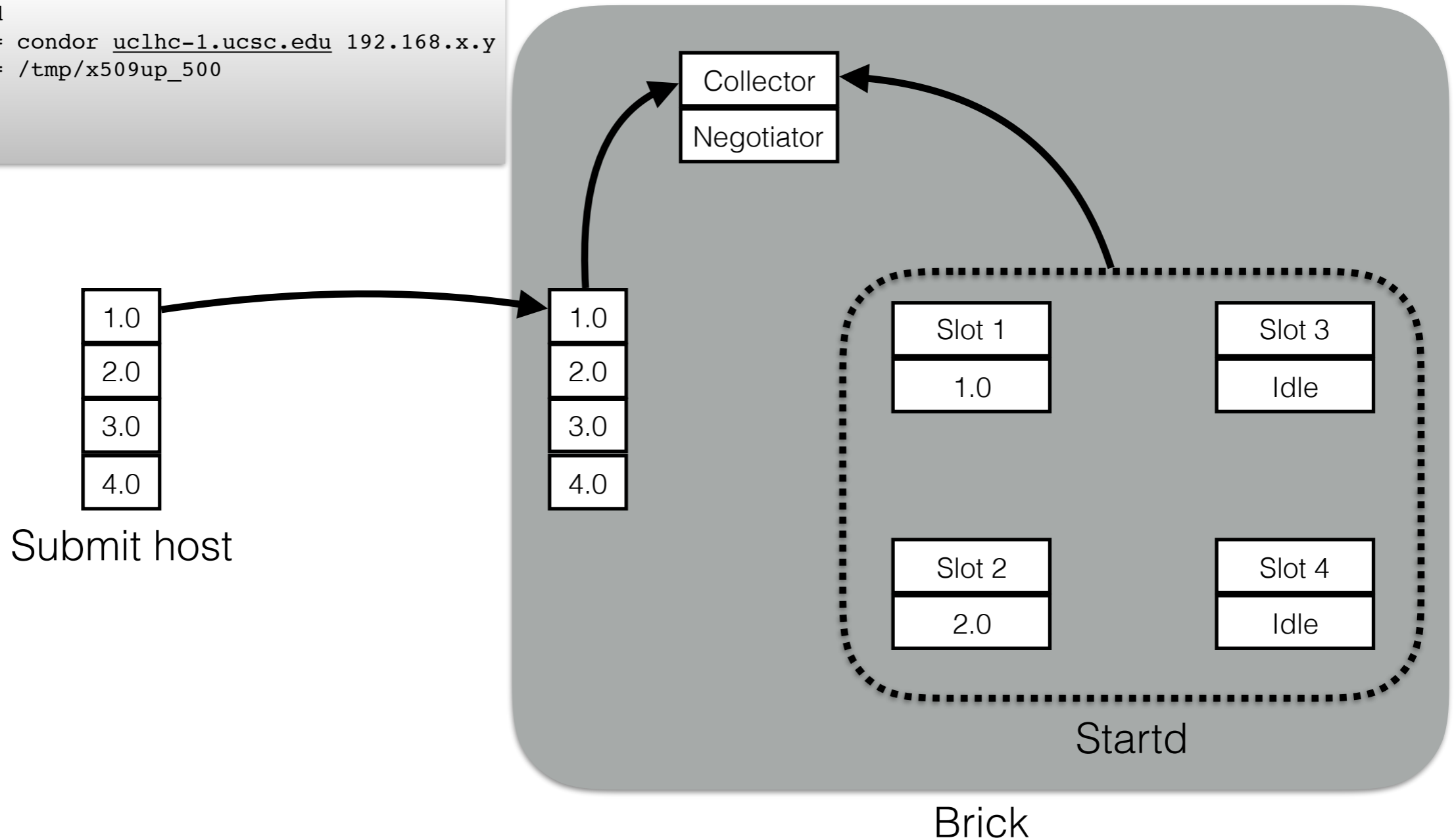
Traditional Submission

universe = vanilla



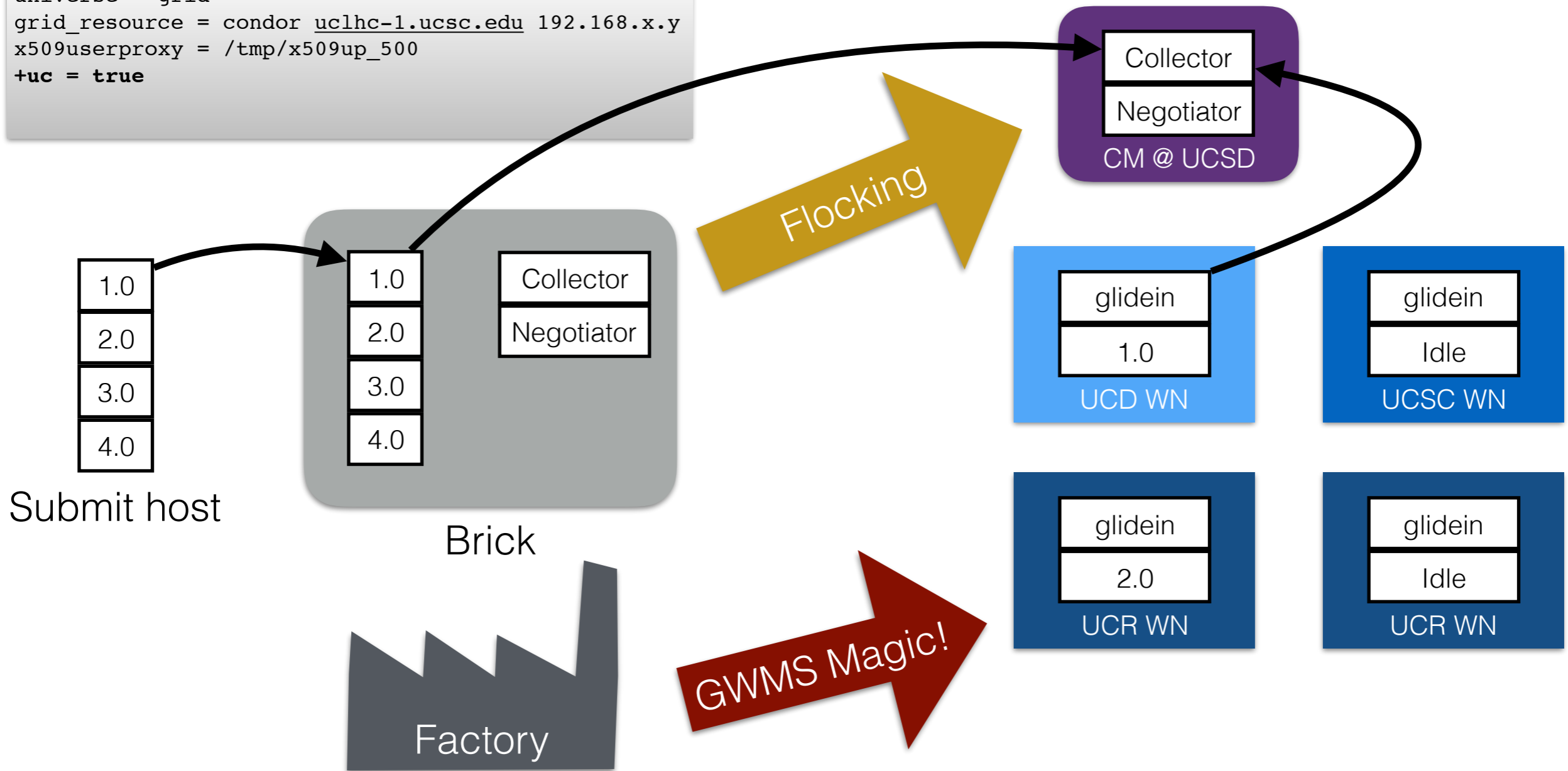
Submit to Brick

```
universe = grid
grid_resource = condor uclhc-1.ucsc.edu 192.168.x.y
x509userproxy = /tmp/x509up_500
+local = true
```



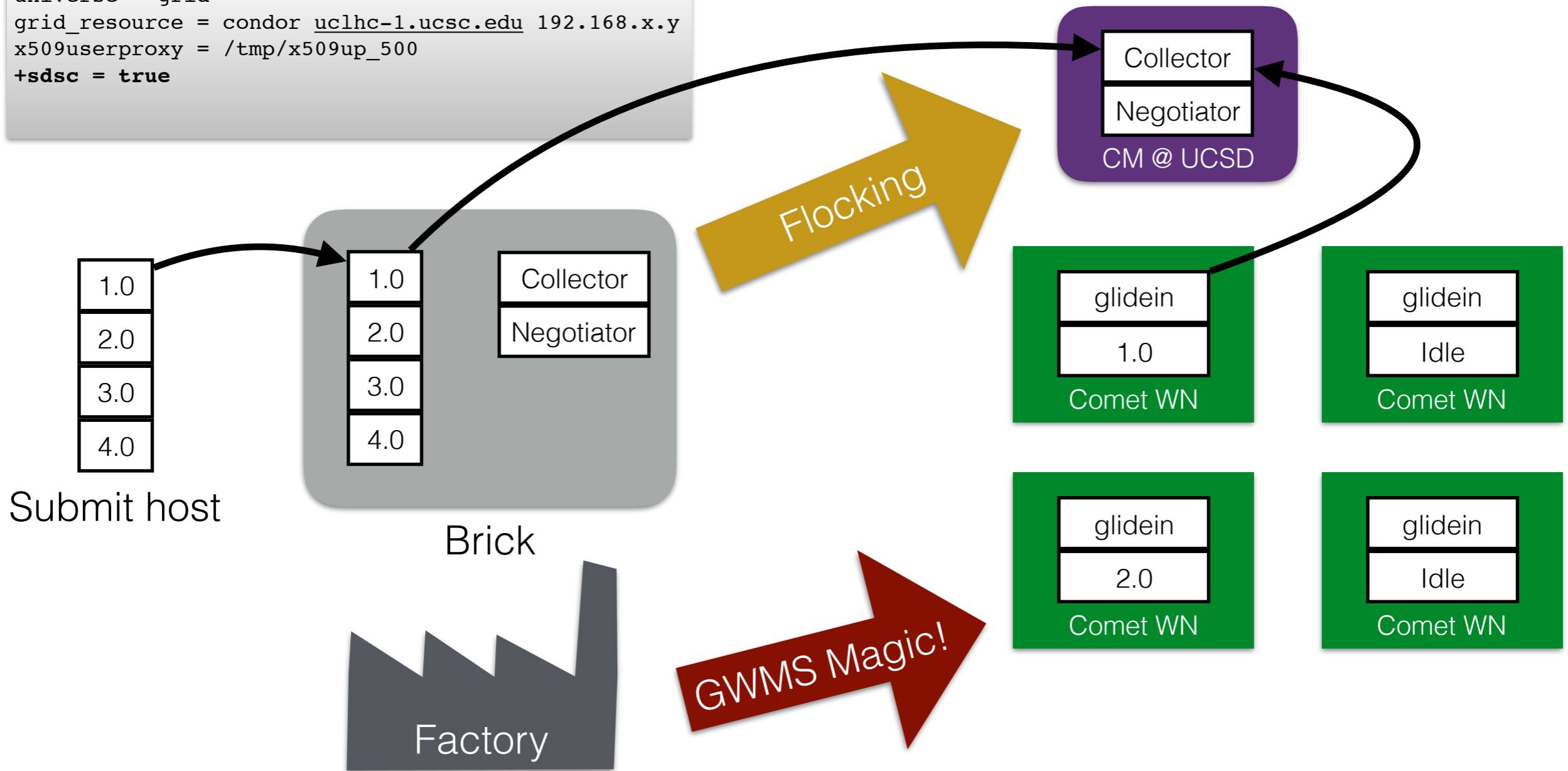
Submit to UCs

```
universe = grid
grid_resource = condor uclhc-1.ucsc.edu 192.168.x.y
x509userproxy = /tmp/x509up_500
+tuc = true
```



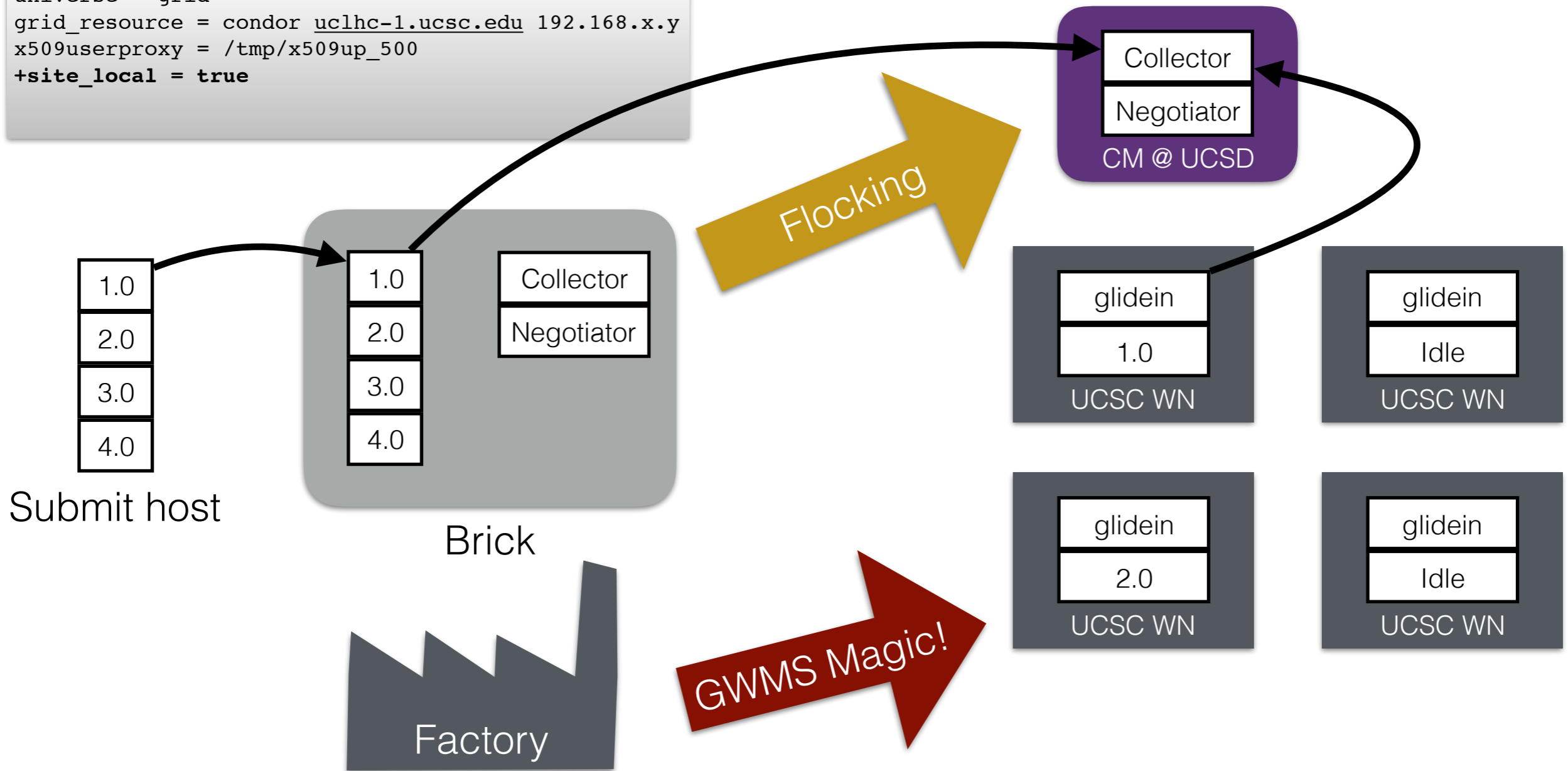
Submit to Comet

```
universe = grid
grid_resource = condor uclhc-1.ucsc.edu 192.168.x.y
x509userproxy = /tmp/x509up_500
+sdsc = true
```



site_local to replace vanilla

```
universe = grid
grid_resource = condor uclhc-1.ucsc.edu 192.168.x.y
x509userproxy = /tmp/x509up_500
+site_local = true
```



Why site_local vs vanilla?

- We dynamically set xrootd cache location using glidein startup scripts:

```
# get cache address based on where we land
glidein_site=`grep -i "^GLIDEIN_Site " $glidein_config | awk '{print $2}'`
CMS_XROOTD_CACHE=`grep -i "$glidein_site"_CMS client/xrootd-cache-location.txt | awk '{print $2}'`
ATLAS_XROOTD_CACHE=`grep -i "$glidein_site"_ATLAS client/xrootd-cache-location.txt | awk '{print $2}'`

# gwms way to set env vars into user job
add_config_line CMS_XROOTD_CACHE "$CMS_XROOTD_CACHE"
add_condor_vars_line CMS_XROOTD_CACHE "S" "-" "+" "Y" "Y" "+"
add_config_line ATLAS_XROOTD_CACHE "$ATLAS_XROOTD_CACHE"
add_condor_vars_line ATLAS_XROOTD_CACHE "S" "-" "+" "Y" "Y" "+"
```

- So users can use env var without having to know where jobs actually land:

```
# Job submitted from UCD lands at UCR:
xrscp root://uclhc-1.ucr.edu:4094//store/mc/RunIIFall15DR76/BulkGravTohhTohVVhbb_narrow_M-900_13TeV-
madgraph/AODSIM/PU25nsData2015v1_76X_mcRun2_asymptotic_v12-v1/10000/40B50F72-5BB4-E511-
A31F-001517FB1B60.root .

# Using env var instead
xrscp root://$CMS_XROOTD_CACHE//store/mc/RunIIFall15DR76/BulkGravTohhTohVVhbb_narrow_M-900_13TeV-
madgraph/AODSIM/PU25nsData2015v1_76X_mcRun2_asymptotic_v12-v1/10000/40B50F72-5BB4-E511-
A31F-001517FB1B60.root .
```

Benefits of Condor-C

- Requires minimal change to existing batch config
 - We provide a drop-in file for `/etc/condor/config.d`
- Users continue using the submit host they are used to logging into
 - We just teach the magic lines to add to the submit file
- Makes user account creation on the brick unnecessary, improving security of DMZ since brick doubles as a data transfer node

Drop-in Condor-C config

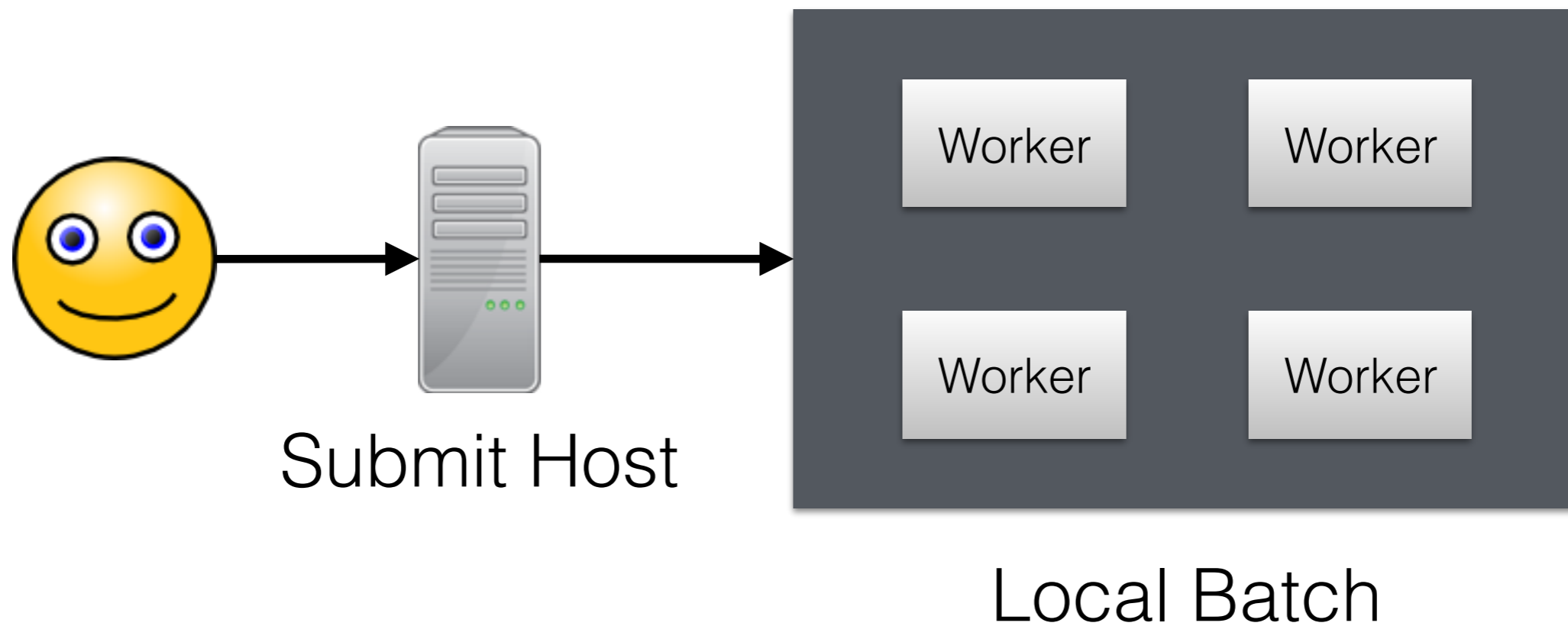
```
# track which site submitted from for accounting and gwms logic (matches GLIDIEN_Site)
SUBMIT_SITE = "UCSC"

# submit attr defaults
local = True
site_local = True
sdsc = False
uc = False
osg = False
SUBMIT_SITE_ATTRS = local site_local sdsc uc osg
SUBMIT_EXPRS = $(SUBMIT_EXPRS) SUBMIT_SITE $(SUBMIT_SITE_ATTRS)

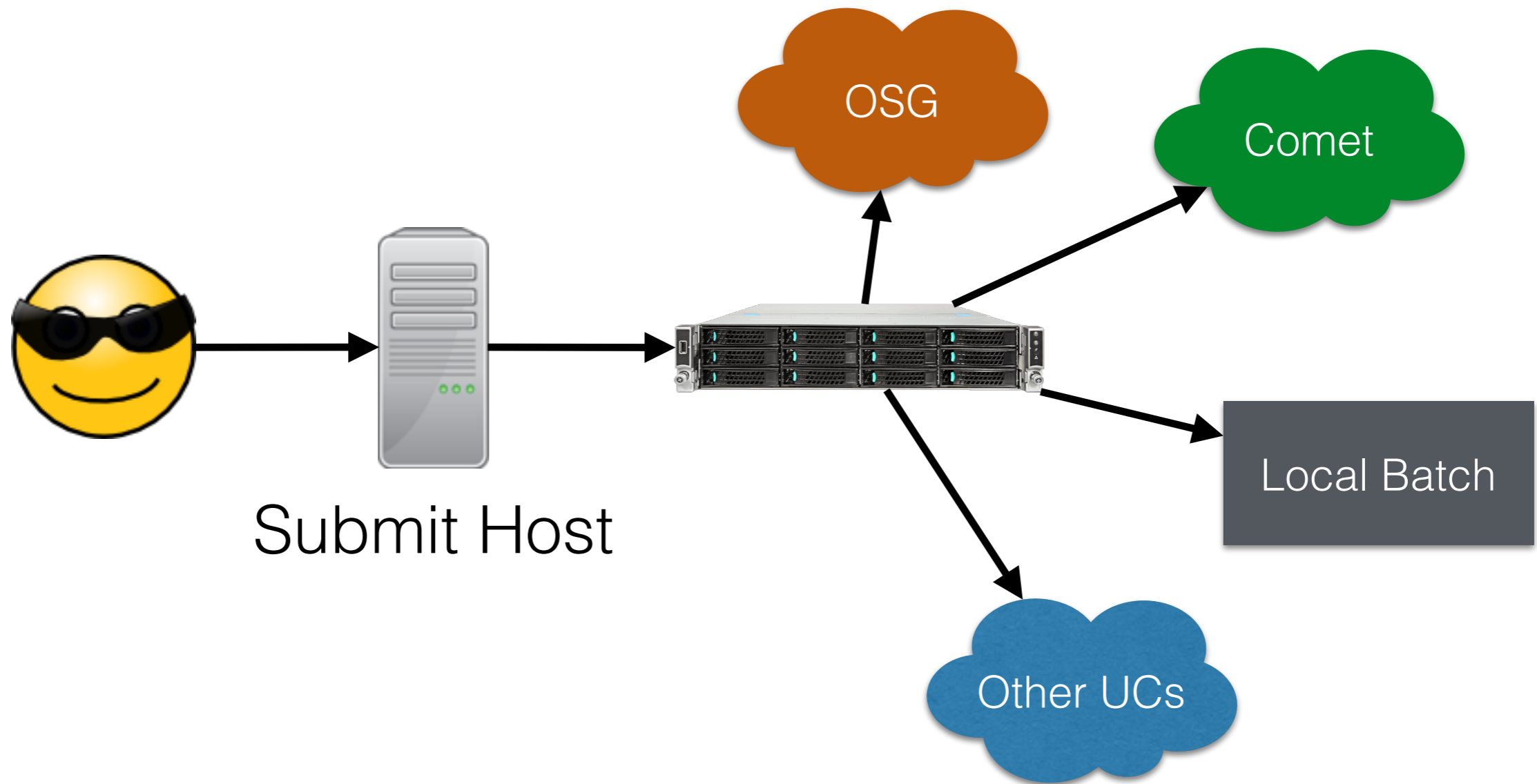
# ensure GSI is enabled for Condor-C
SEC_CLIENT_AUTHENTICATION_METHODS = GSI, $(SEC_CLIENT_AUTHENTICATION_METHODS)

# don't let GSI complain if submitting via brick local network IP
GSI_SKIP_HOST_CHECK_CERT_REGEX = ^\//DC\=org\//DC\=opensciencegrid\//O\=Open\ Science\ Grid\//OU\=Services
\//CN\=uclhc\-1\.ucsc\.edu$
```

Traditional T3

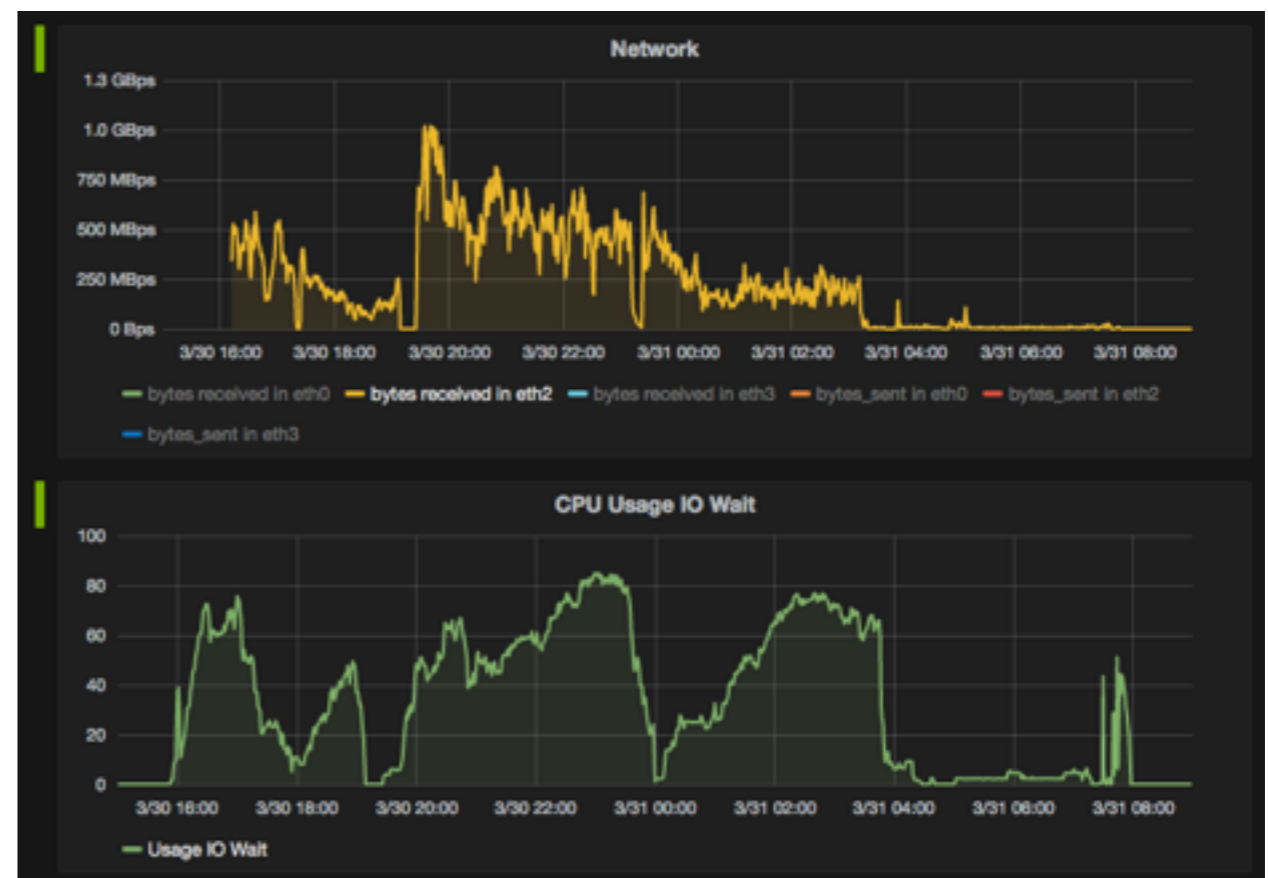
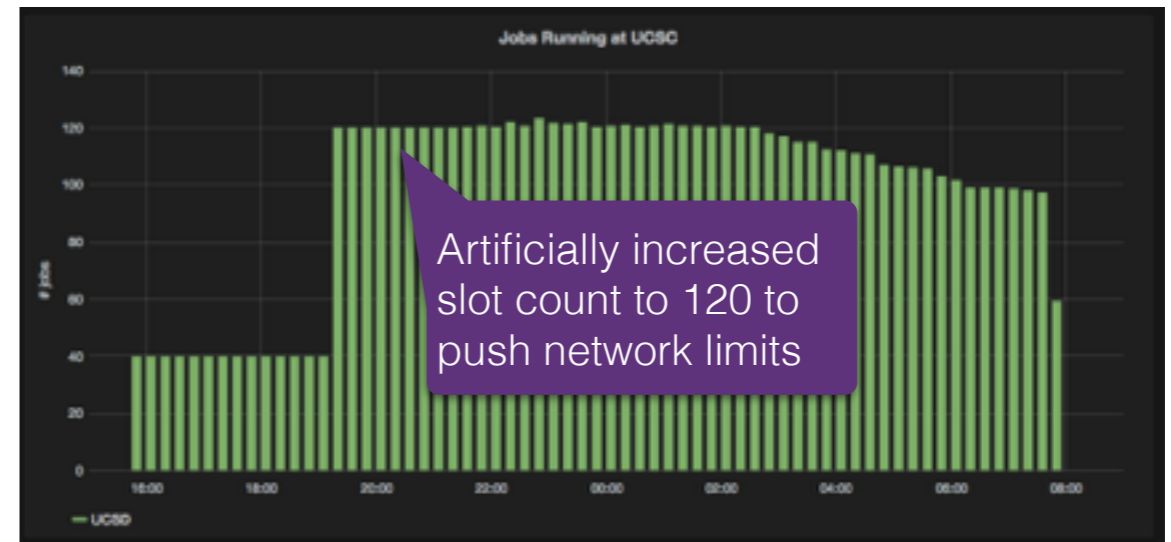


Uber Tier 3



UCSC Cache Stress Test

- Test jobs running xrdcp bring data in from CMS AAA data federation
- As expected initially we saturated the inbound at 10 gbps
- Over time network usage reduces as we access already fetched files from the disk cache
- Network choppiness and disk thrashing was caused by a subset of files xrootd was failing to fetch
 - O(1000) lingering TCP connections observed between cache and data servers long after clients disconnected from cache
 - We suspect a bug in the xrootd proxy cache when accessing broken files, more testing is needed



Attn HTCondor Devs!

(aka Edgar's wish list)

- Surprises but not critical (Condor-C):
 - Unable to get Condor-C to use password auth
 - Started thread with HTCondor support
 - Not a big issue since our users already have grid certs, so we chose GSI
 - Condor-C doesn't use local network unless you give local hostname / ip in submit file
 - SUBMIT_EXPRS has to be set on submission side schedd
 - variable is ignored on remote schedd side
- Improvements we definitely want:
 - Very interested in the generalization of monitoring metrics, we dropped Ganglia in favor of **Influx DB**
 - Currently guilty of periodic condor-q's to parse running user job numbers (sorry Brian!)
- An "it would be nice":
 - Any plans on implementing expression based flocking?

Conclusion

- LHC @ UC project utilizes the PRP network to enhance the T3s at each site by providing:
 - A unified way to submit locally and compute globally
 - The ability to decouple data placement from where the jobs run
- The central management of the services by dedicated admins at UCSD allows the local UC users to worry less about infrastructure maintenance and focus more on getting science done

One HTCondor pool to rule them all



Questions?