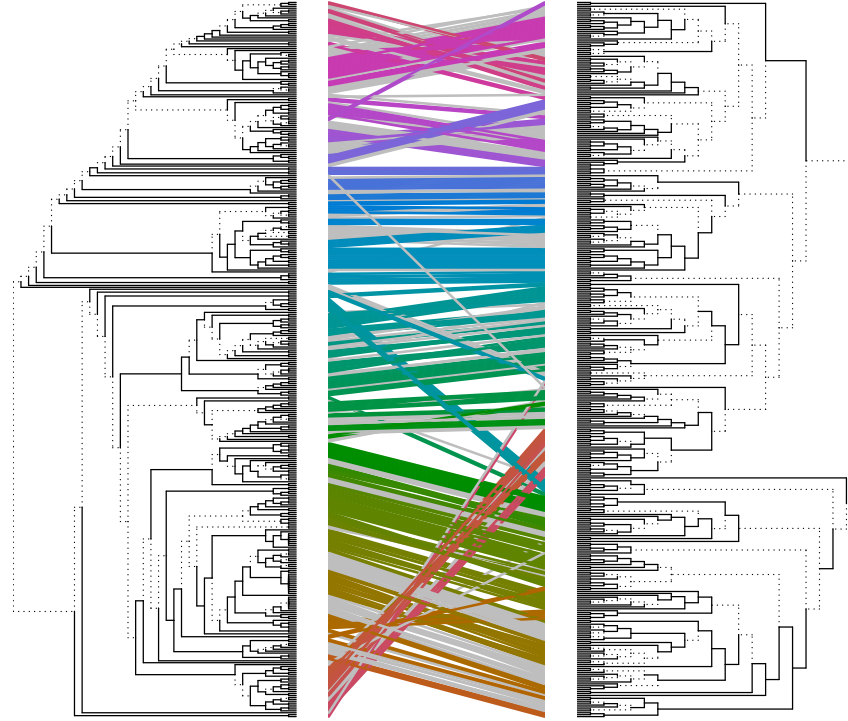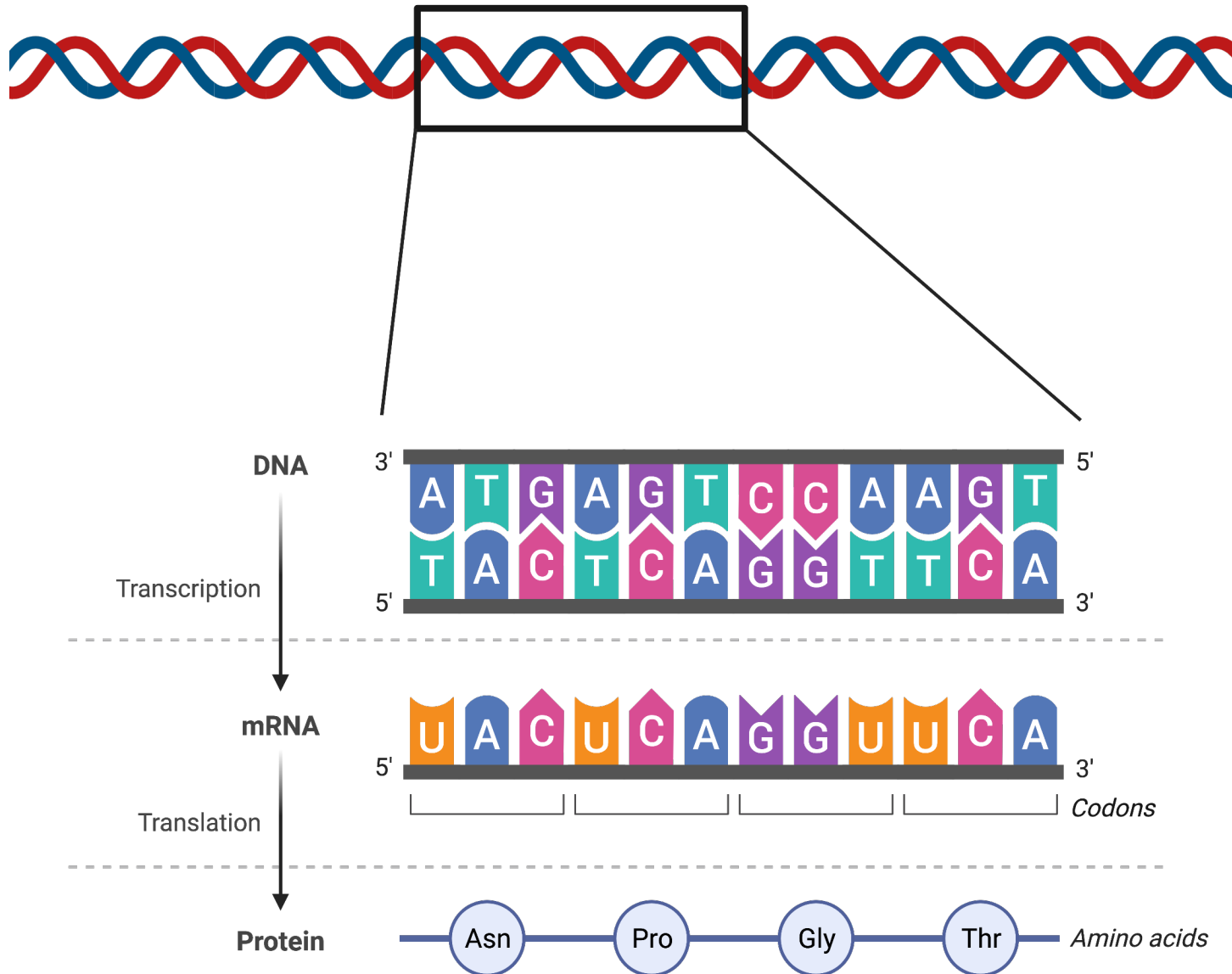# Computational Biology on the Open Science Grid

Nicholas Cooley
Wright Lab
University of Pittsburgh
Department of Biomedical Informatics

It is relatively straightforward to predict where genes are in genomes, even newly sequenced, novel isolates.
Figuring out what new genes do however, is not trivial.

Created with BioRender.com

Classification / Annotation:
Does this sequence have the same function or job as a sequence in some training data whose function or job is known?

Give a novel sequence a descriptive and succinct label that represents that sequence's function.

>LibrarySequence01
MQRNRLFSENTTELMSTPHHD...

>LibrarySequence02
MAIRQWMMIGKHLCRFELRRF...

>LibrarySequence03
MHLWPWIMQDEFEVAMCWRQK...

>LibrarySequence04
MSQWPSNERMEANDDGRTGYS...

>LibrarySequence05
MKIHKLTPCEFMENRSQYKYA...

>LibrarySequence06
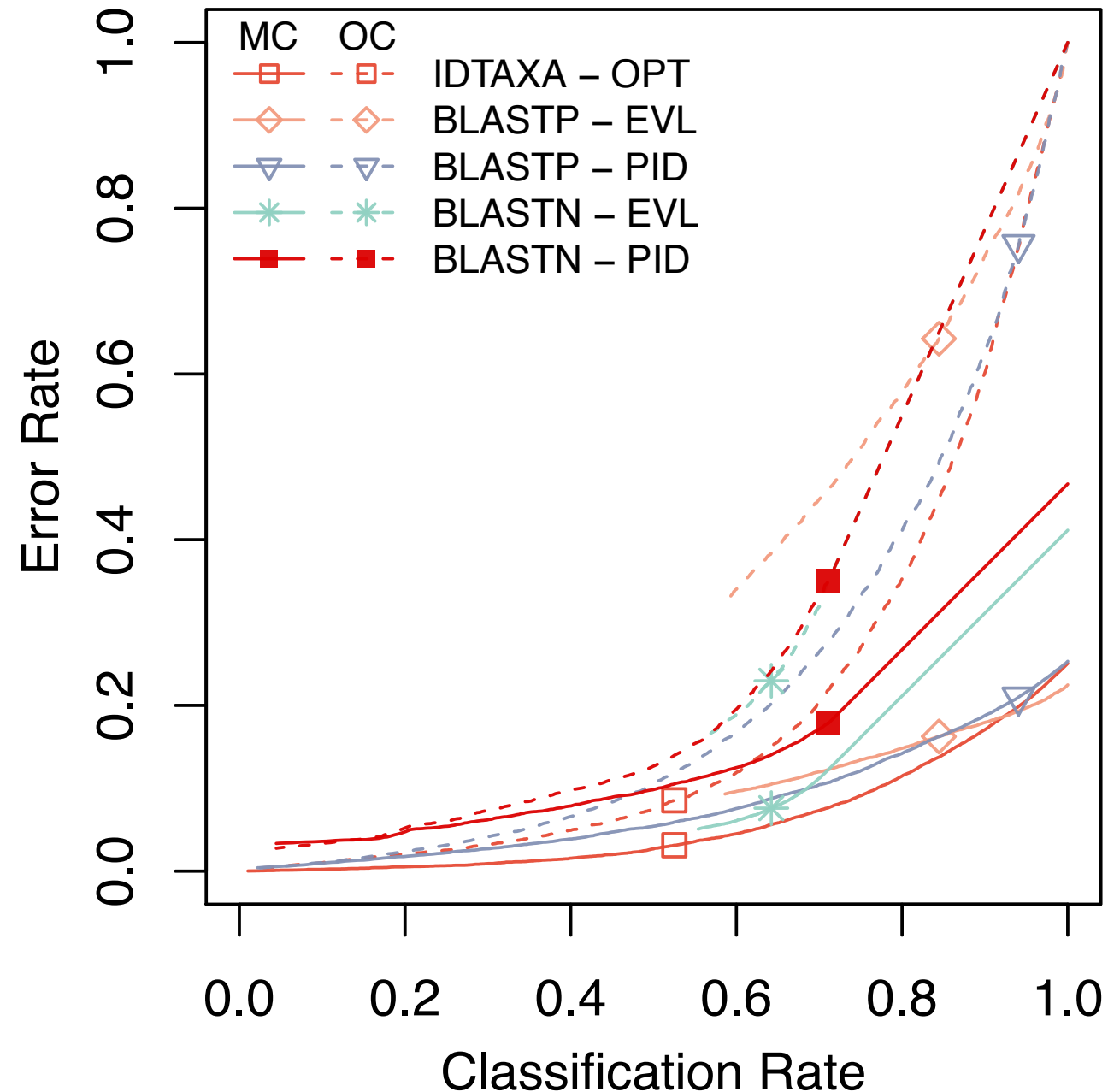MDKKWYYKWEMRQECDPRSVD...

>LibrarySequence07
MNCWHTWMMKDRRNIGETCHM...

>LibrarySequence08
MFRARYHMPHTCYESGPMHKD...

>NewSequence01
MSADDHGMRNVPKHIFNKGLK...

Does **this** sequence have a
function with a representative
in **this** library?

TL;DR we built a classifier:
- Accurate functional classification is difficult
- Emphasis on conservative classification to avoid overclassification of truly novel sequences
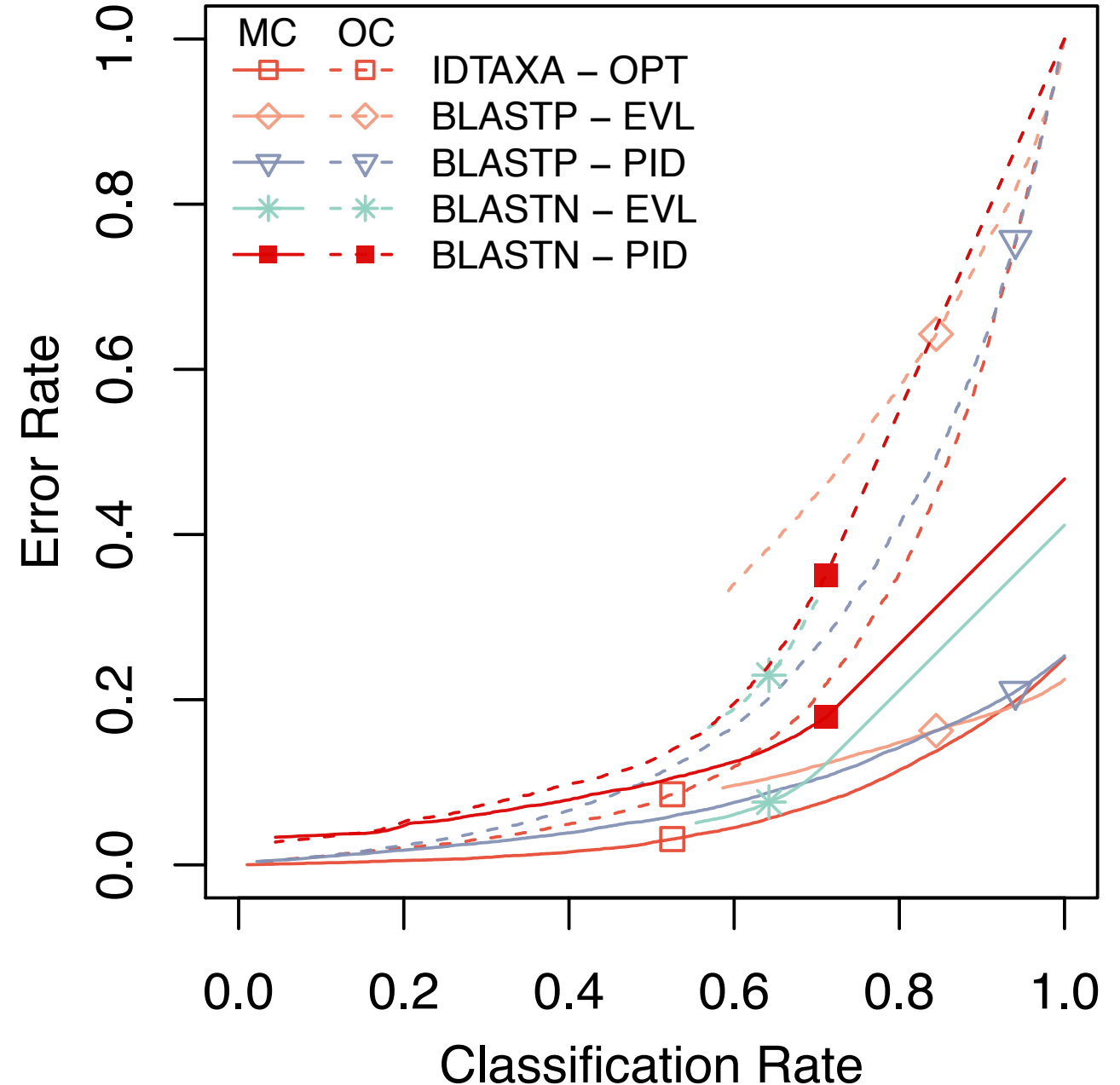
Publication coming soon!

This work would not have been
possible without the open science
grid:
• Data preparation
• Parameter tuning
  • k-mer characteristics
• Cross validation
• Testing, testing, testing …

Publication coming soon!

One last bit of biology to introduce:

>Sequence01
MSAD DHGMRNVPKHIFNK GLKHWPKYRPITWQLSDFGEWEFDS

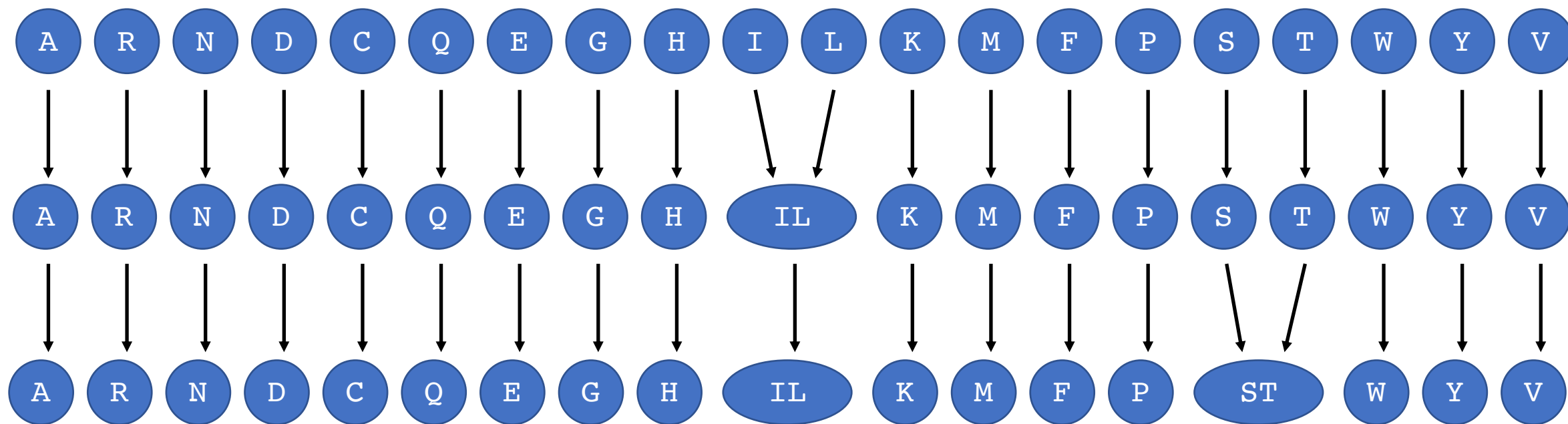If **this** word/k-mer appears **here** there is implied significance.

What about with minor changes?
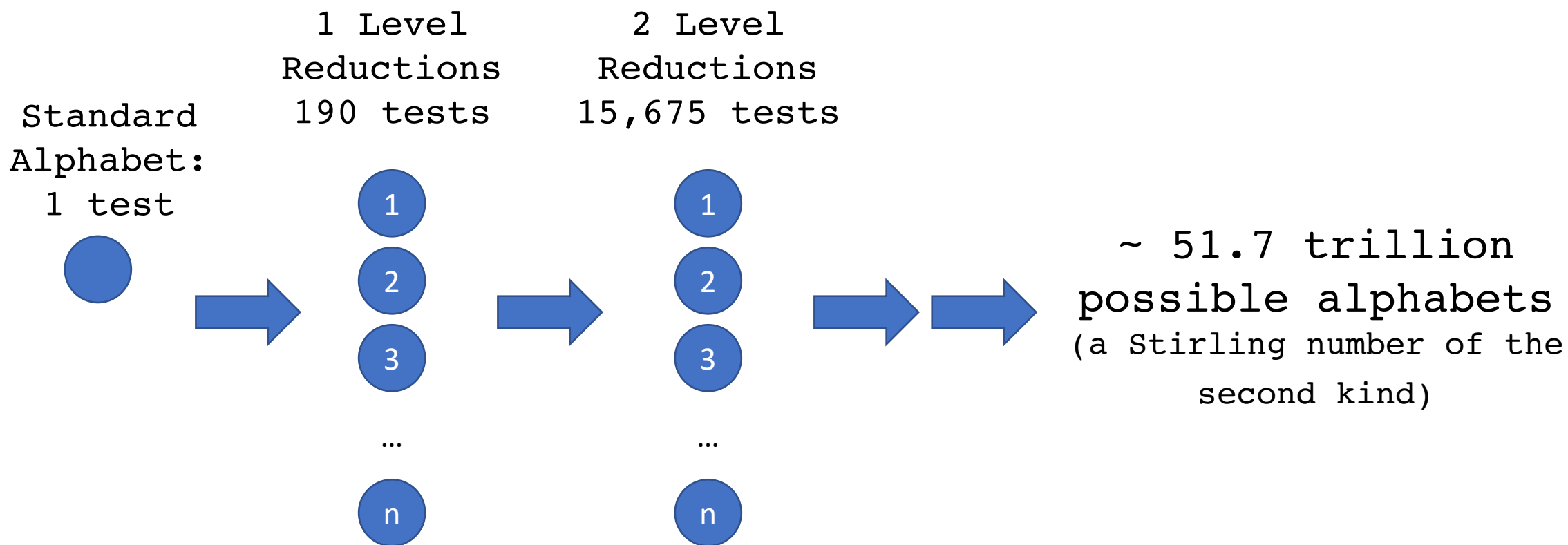
What about with major changes?

>Sequence02
MDQKMGDQCTP DHGMRNVPKHIFNK YPASTNEKDHYNMLDGAVNE

Performance of our classifier improves when the standard amino acid alphabet is substituted with a reduced alphabet:
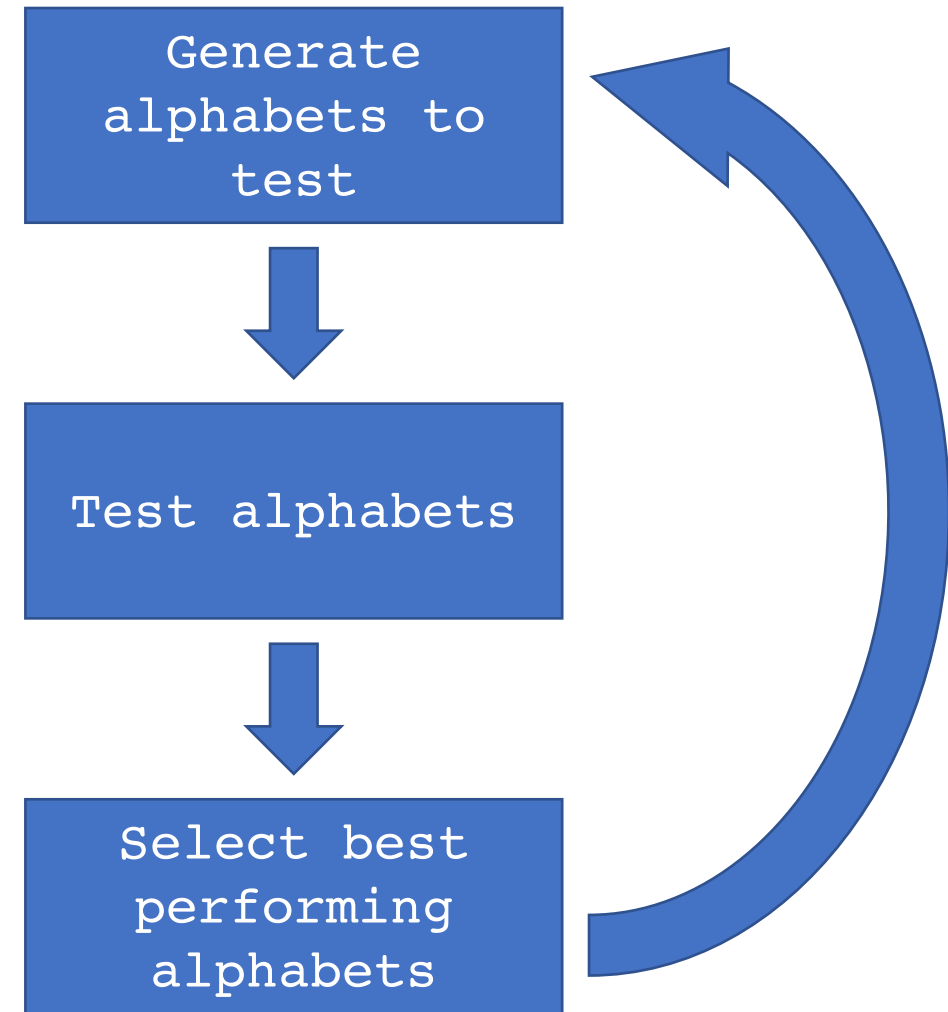


Test alphabet performance, perform a reduction, test again, repeat ad nauseum.

# If reduced alphabets provide improved performance, how do we select the *best* reduced alphabet?

Standard
Alphabet:
1 test

1 Level
Reductions
190 tests

1
2
3
…
n

2 Level
Reductions
15,675 tests

1
2
3
…
n

~ 51.7 trillion
possible alphabets
(a Stirling number of the
second kind)

51.7 trillion tests is probably
too many tests.
- Iterate down through alphabet
  sizes
- Only test reductions of highest
  performing alphabets from
  previous level
- Avoid brute force testing of
  every possible alphabet

```
┌──────────────────┐
│    Generate       │
│ alphabets to     │
│    test           │
└──────────────────┘
         │
         ▼
┌──────────────────┐
│  Test alphabets   │
└──────────────────┘
         │
         ▼
┌──────────────────┐
│   Select best     │
│   performing      │
│   alphabets       │
└──────────────────┘
```
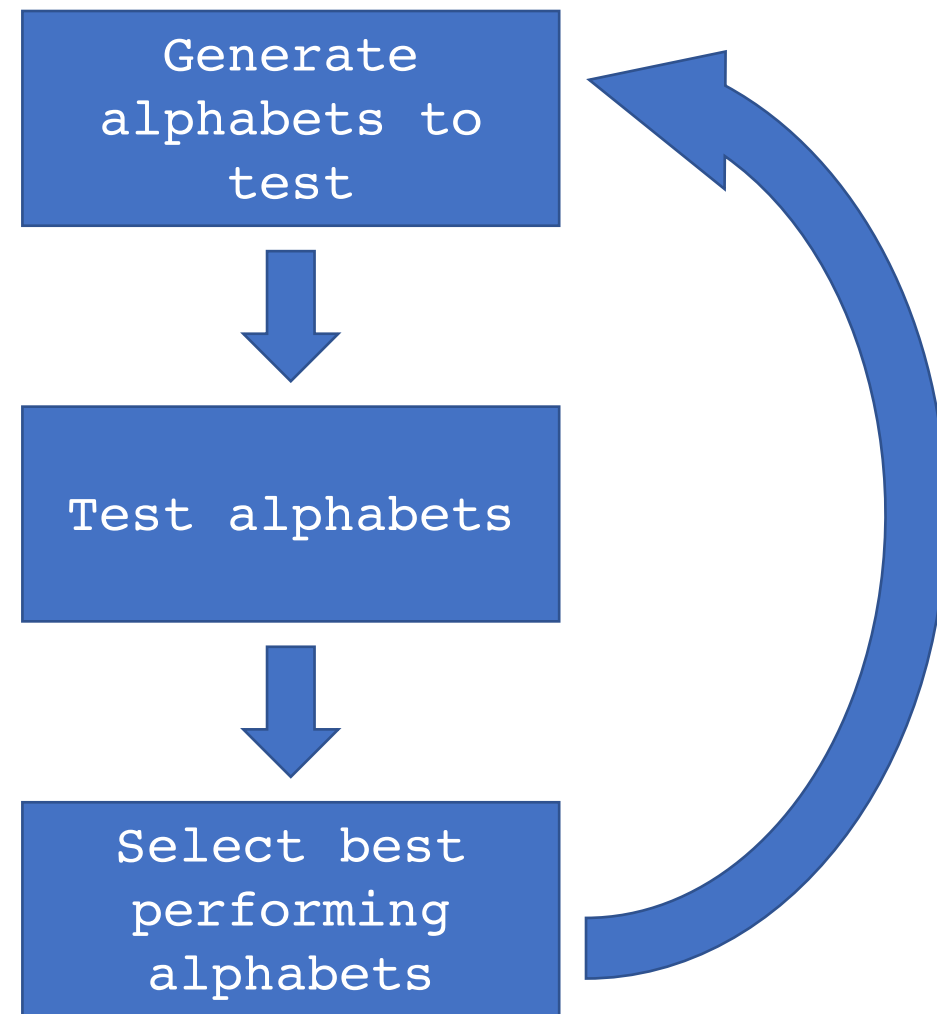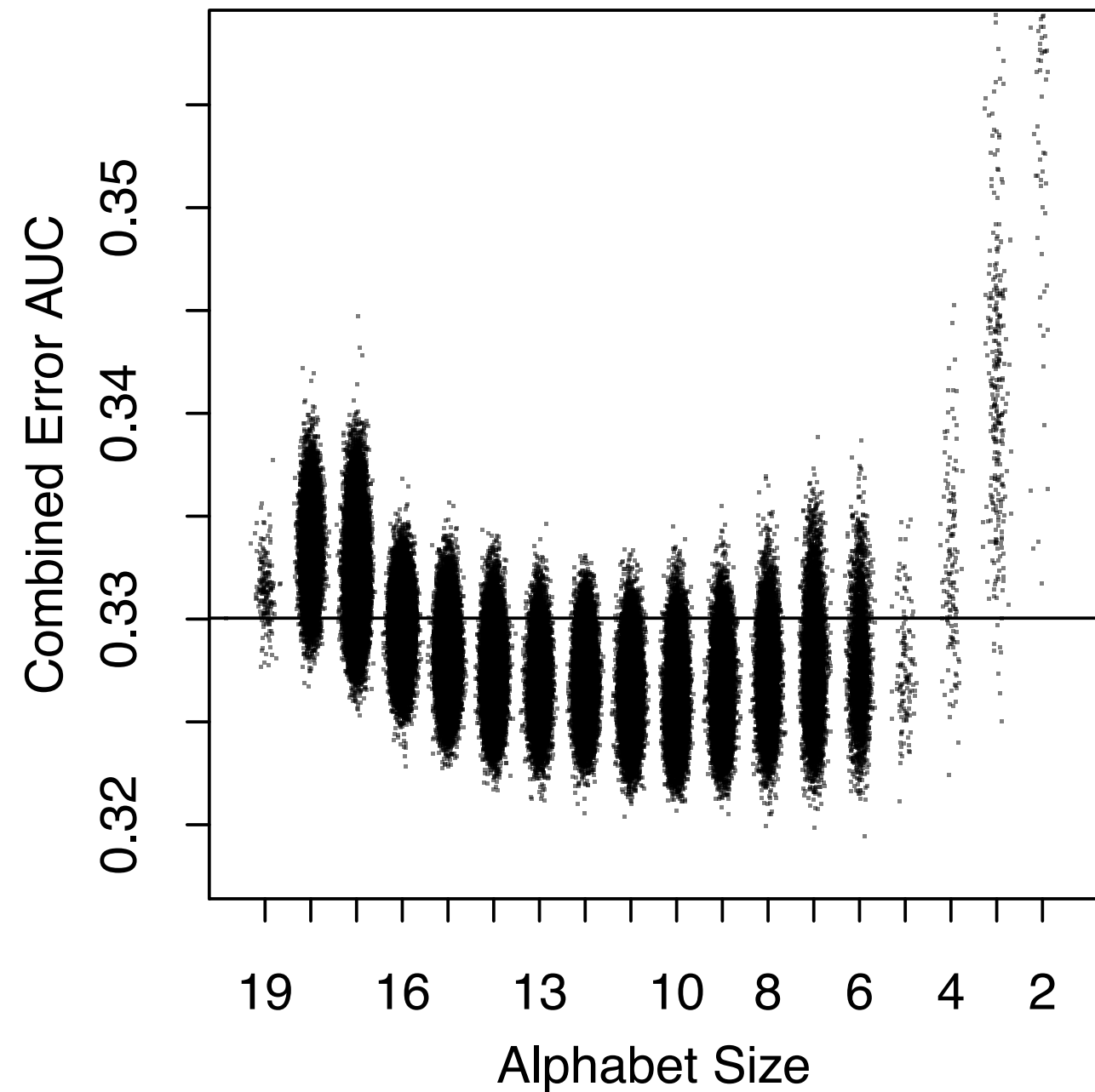
# We've got a DAG for that …

Test, Consolidate, Repeat
- Testing alphabets has modest requirements:
  - 1 CPU
  - 1 GB disk
  - > 4 GB memory
- Consolidating results at each level has trivial requirements:
  - 1 CPU
  - 1 GB disk
  - 2 GB memory
    - During consolidation, parameters for next level are set

```
1   JOB A OSG01Job.sh
2   JOB B OSG01Consolidate.sh
3   JOB C OSG02Job.sh
4   JOB D OSG02Consolidate.sh
5   JOB E OSG03Job.sh
6   JOB F OSG03Consolidate.sh
7   JOB G OSG04Job.sh
8   JOB H OSG04Consolidate.sh
9   JOB I OSG05Job.sh
10  JOB J OSG05Consolidate.sh
11  JOB K OSG06Job.sh
12  JOB L OSG06Consolidate.sh
13  JOB M OSG07Job.sh
14  JOB N OSG07Consolidate.sh
15  JOB O OSG08Job.sh
16  JOB P OSG08Consolidate.sh
17  JOB Q OSG09Job.sh
18  JOB R OSG09Consolidate.sh
```

# We've got a DAG for that …
# And we try to keep it simple:

```bash
1   #!/bin/bash
2
3   # running inside a singularity container, ENV commands in the dockefile aren't r
4   # This path needs to match where the executable was installed in the dockerfile
5   # export PATH=/blast/ncbi-blast-2.9.0+/bin:$PATH
6
7   Rscript JobScript.R $1 Level01Alphabets.RData
8
9   if [ -e Result*.RData ]
10  then
11      exit 0
12  else
13      exit 1
14  fi
15  
```
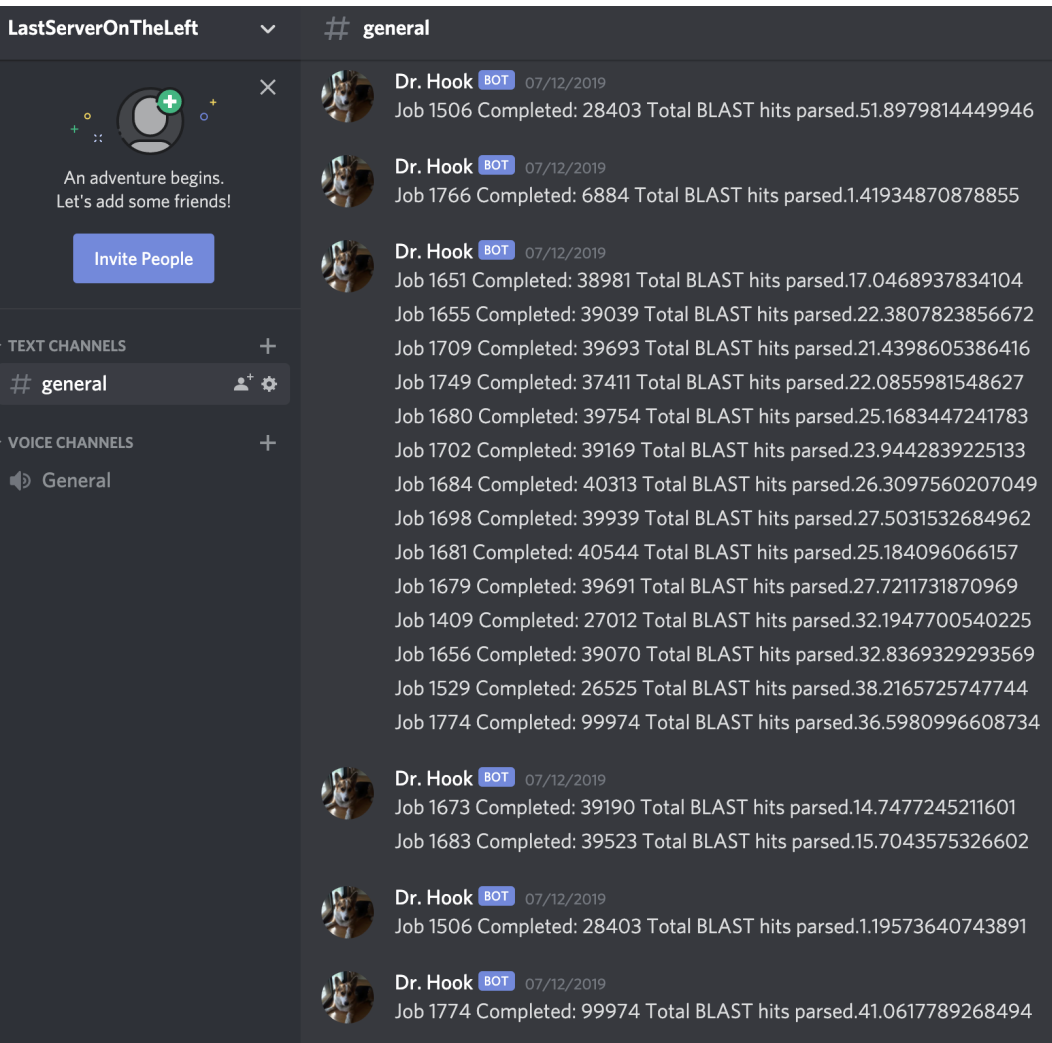
# Other cool tidbits:

|     | 1 | 2 | 3 | 4 | 5 | 6 | … | n |
|-----|---|---|---|---|---|---|---|---|
| **1** | | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| **2** | | | ■ | ■ | ■ | ■ | ■ | ■ |
| **3** | | | | ■ | ■ | ■ | ■ | ■ |
| **4** | | | | | ■ | ■ | ■ | ■ |
| **5** | | | | | | ■ | ■ | ■ |
| **6** | | | | | | | ■ | ■ |
| **…** | | | | | | | | ■ |
| **n** | | | | | | | | |

A typical job set up for us: Perform all pairwise comparisons in a set of genomes — give 1 comparison to each node.

With relatively trivial requirements for nodes (1 GB disk, 2 GB memory, 1 CPU) we can complete ~ 70,000 jobs at 10 minutes per job in a weekend.

# Other cool tidbits:



Monitoring jobs in real time is complicated.

`npcooley$ watch —n 5 condor_q`

But we're not always at a work computer and ssh'd into our login node

Discord can collect results for us but …

# Acknowledgements



Wright Lab:
- Erik Wright
- Maria Bond
- Andrew Beckley
- Allison Petrick
- Sam Blechman
- Shania Khatri
- Nishant Panicker

Open Science Grid:
- OSG User School
- Christina Koch
- Lauren Michael
- Mats Rynge
- Carrie Brown

1) Pordes, R. et al. (2007). "The Open Science Grid", J. Phys. Conf. Ser. 78, 012057.doi:10.1088/1742-6596/78/1/012057.
2) Sfiligoi, I., Bradley, D. C., Holzman, B., Mhashilkar, P., Padhi, S. and Wurthwein, F. (2009). "The Pilot Way to Grid Resources Using glideinWMS", 2009 WRI World Congress on Computer Science and Information Engineering, Vol. 2, pp. 428–432. doi:10.1109/CSIE.2009.950.