# Performance and impact of dynamic data placement in ATLAS

Thomas Maier[1]
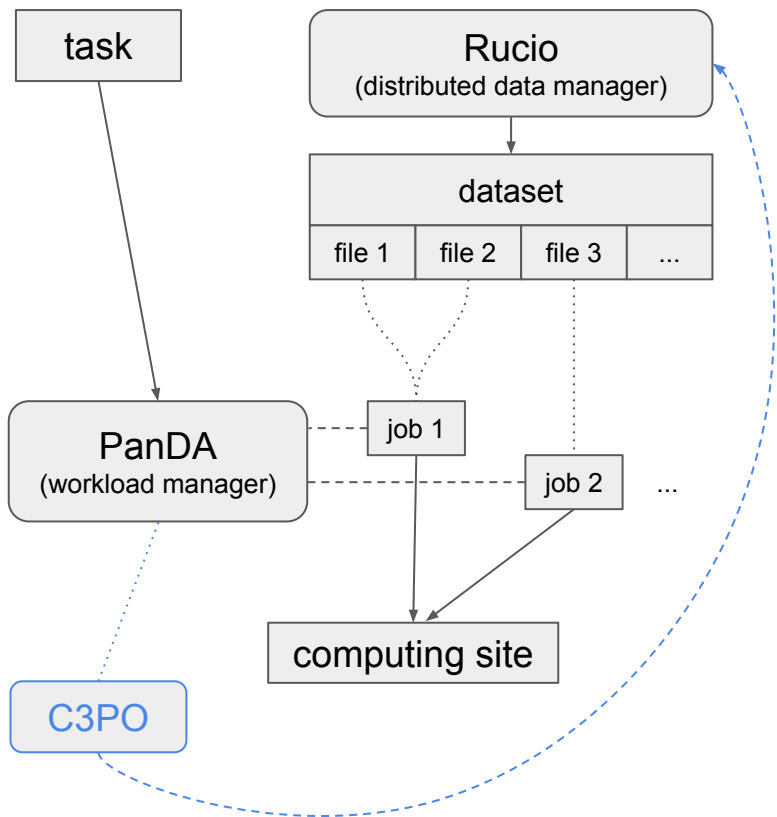on behalf of the ATLAS Collaboration

Co-authors: T. Beermann[2], G. Duckeck[1], M. Lassnig[3], F. Legger[1], M. Magoni[4], I. Vukotic[5]

[1]Ludwig-Maximilians-Universität München, [2]Universität Innsbruck, [3]CERN, [4]Politecnico di Torino, [5]University of Chicago
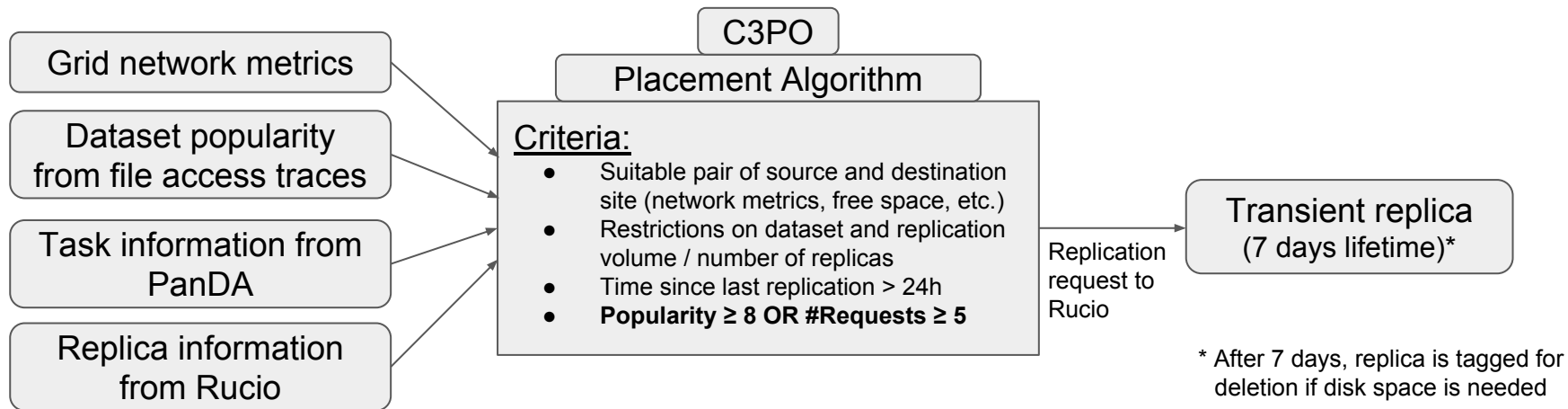
CHEP 2018 Conference - 11.07.2018

# Introduction



- Why dynamic data placement?
  - Input data popularity can vary heavily dependent on data type/format or time period
  - By default, numbers of dataset replicas are statically set ⇒ a given dataset is only available at a couple of computing sites
  - Temporary high demand of datasets can lead to delayed processing of jobs ⇒ high waiting times for physics analyses
- In ATLAS, new **dynamic data placement agent C3PO** was developed and operates during Run-2 → CHEP 2016 contribution
- Analysis of C3PO performance and impact using the ATLAS/CERN analytics infrastructure (ElasticSearch, Kibana, Hadoop, etc.) → dedicated talk on ATLAS analytics platforms

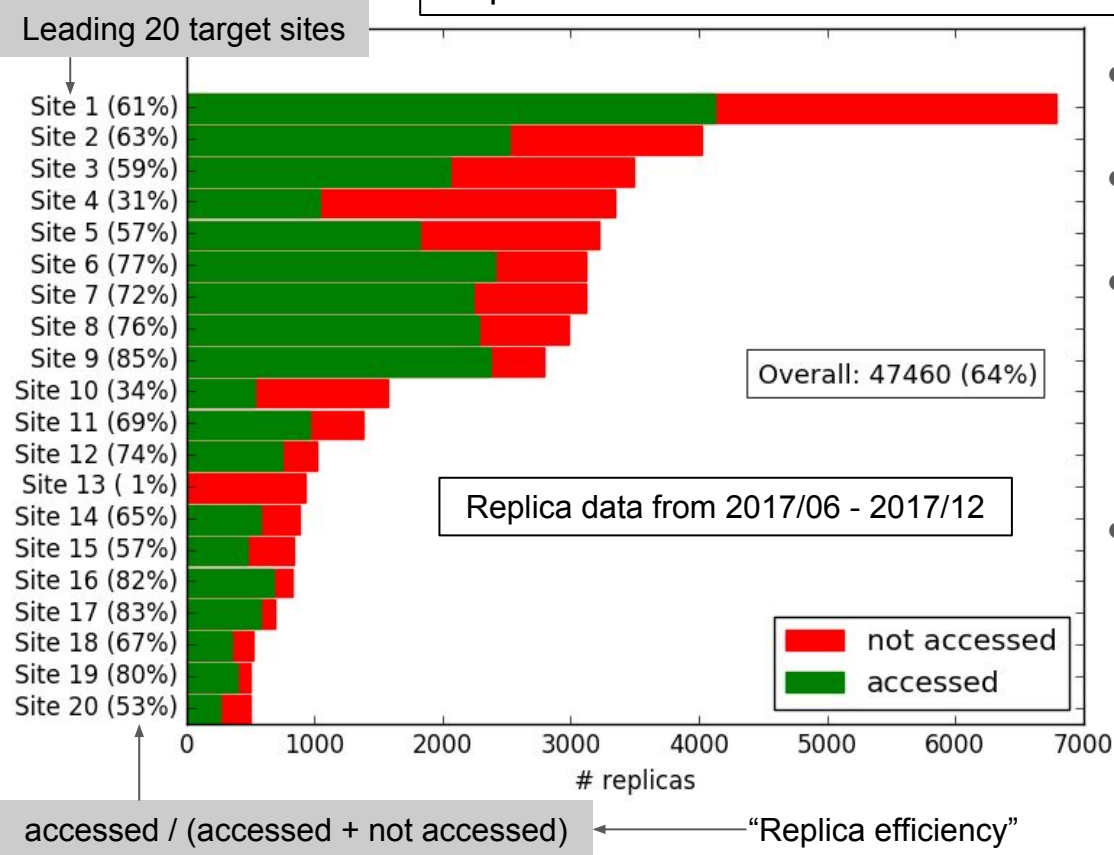# C3PO - Dynamic data placement agent in ATLAS

| Grid network metrics |
| --- |

| Dataset popularity from file access traces |
| --- |

| Task information from PanDA |
| --- |

| Replica information from Rucio |
| --- |

**C3PO**

**Placement Algorithm**

<u>Criteria:</u>
- Suitable pair of source and destination site (network metrics, free space, etc.)
- Restrictions on dataset and replication volume / number of replicas
- Time since last replication > 24h
- **Popularity ≥ 8 OR #Requests ≥ 5**

Replication request to Rucio

| Transient replica (7 days lifetime)* |
| --- |

\* After 7 days, replica is tagged for deletion if disk space is needed

- <u>Popularity:</u> number of dataset accesses in the past 7 days
- <u>#Requests:</u> number of user analysis tasks that use the dataset as input, submitted in the past 24h

# Replica usage after creation - Sites

Replicas which were **accessed** or **not accessed** after creation

Leading 20 target sites



Site 1 (61%)
Site 2 (63%)
Site 3 (59%)
Site 4 (31%)
Site 5 (57%)
Site 6 (77%)
Site 7 (72%)
Site 8 (76%)
Site 9 (85%)
Site 10 (34%)
Site 11 (69%)
Site 12 (74%)
Site 13 ( 1%)
Site 14 (65%)
Site 15 (57%)
Site 16 (82%)
Site 17 (83%)
Site 18 (67%)
Site 19 (80%)
Site 20 (53%)

Overall: 47460 (64%)

Replica data from 2017/06 - 2017/12

not accessed
accessed

# replicas

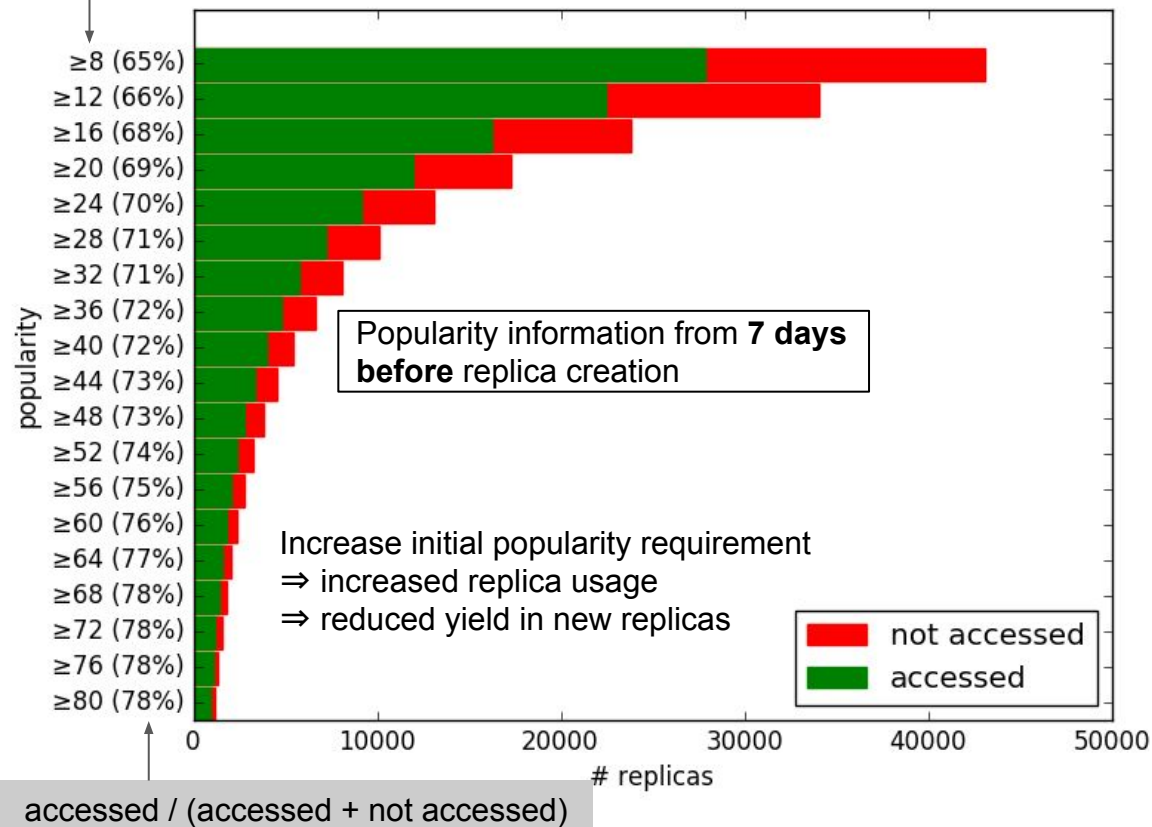accessed / (accessed + not accessed) ← "Replica efficiency"

- File access information from Rucio traces ⇒ replicas were used or not?
- Overall, 64% of replicas were used after their creation
- Target site selection priority is weighted to avoid uneven distribution of replicas
  - Still clustering at a few sites
  - Possibly periods of high disk space availability
- Replica efficiency can strongly depend on where it was put
  - Placement algorithm doesn't take into account computing resources at target site
  - Some correlations with data type or format can be seen as well
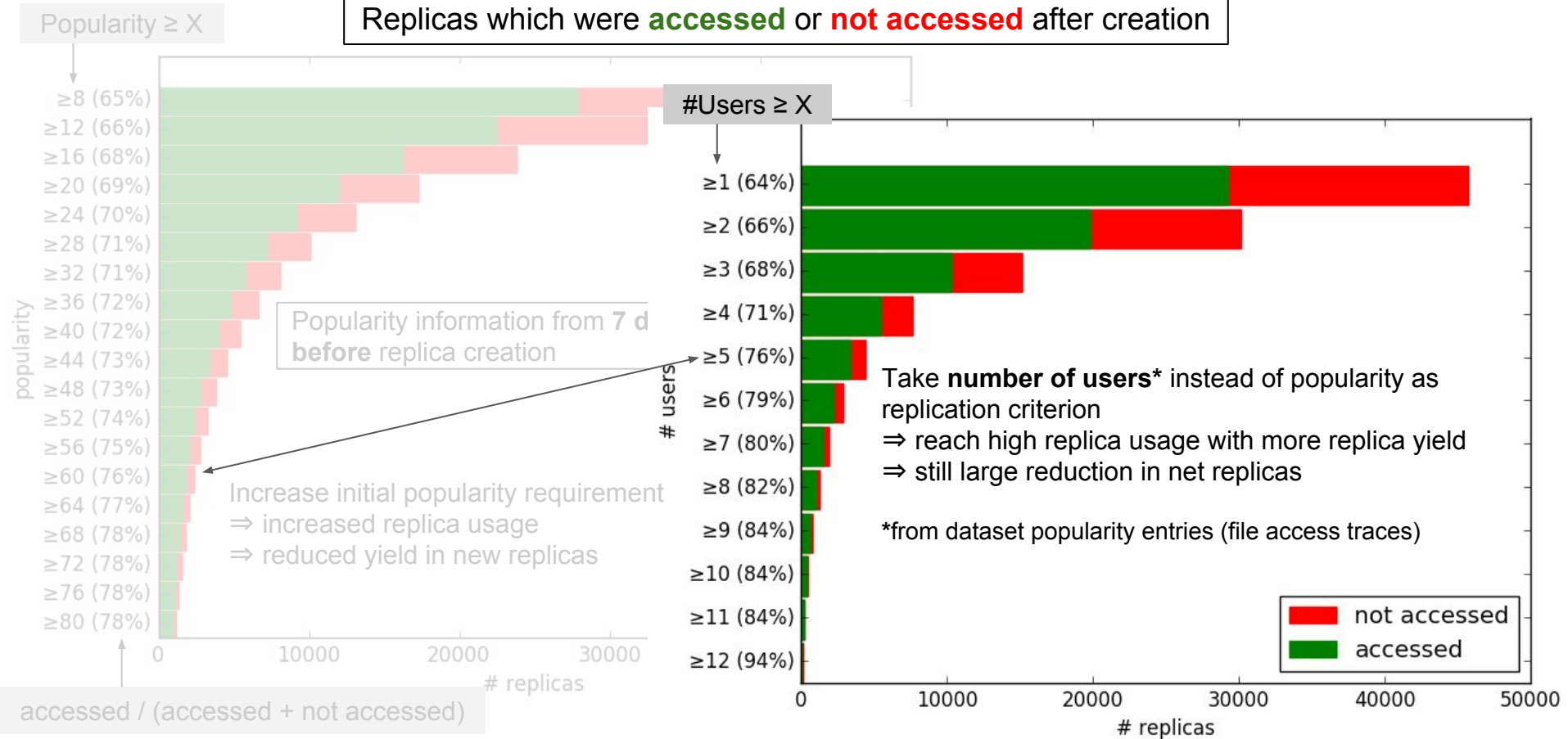
# Replica usage after creation - Popularity



Popularity ≥ X

Replicas which were **accessed** or **not accessed** after creation

Popularity information from **7 days before** replica creation
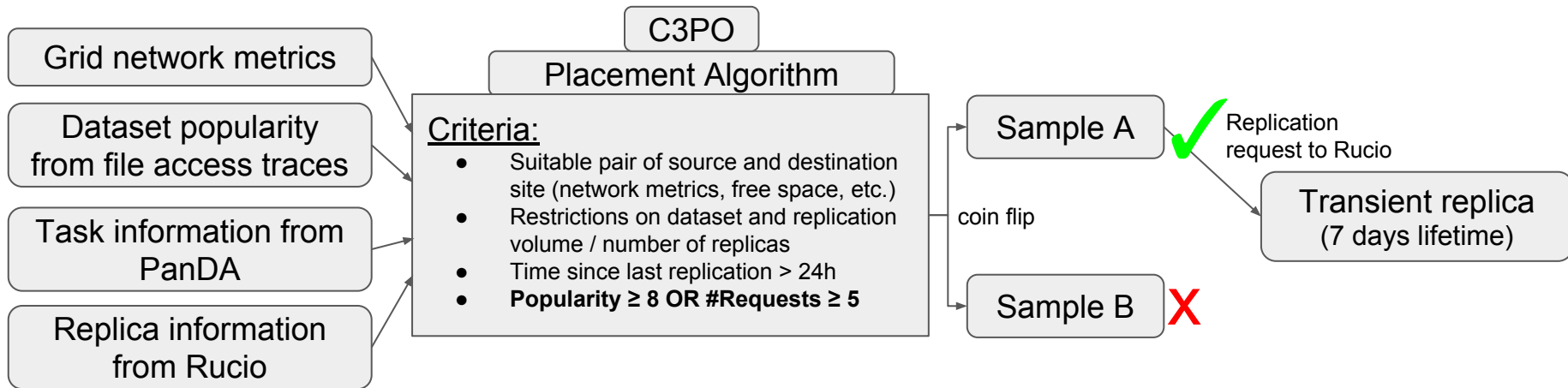
Increase initial popularity requirement
⇒ increased replica usage
⇒ reduced yield in new replicas

accessed / (accessed + not accessed)

# Replica usage after creation - Popularity

Replicas which were **accessed** or **not accessed** after creation

Popularity ≥ X

#Users ≥ X

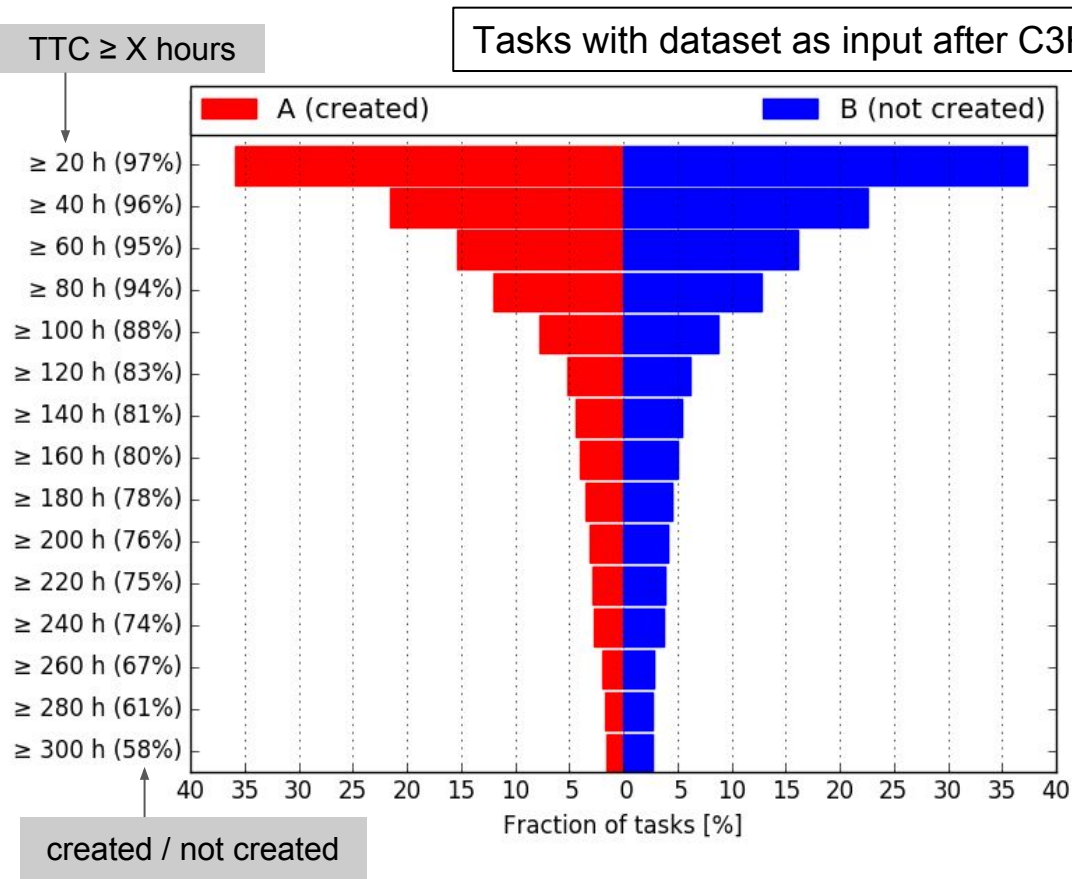Popularity information from **7 d before** replica creation

Increase initial popularity requirement
⇒ increased replica usage
⇒ reduced yield in new replicas

accessed / (accessed + not accessed)

Take **number of users\*** instead of popularity as replication criterion
⇒ reach high replica usage with more replica yield
⇒ still large reduction in net replicas

\*from dataset popularity entries (file access traces)

# C3PO impact analysis

- Attempt to measure effect of C3PO operations on Grid throughput in terms of processed tasks and jobs
- Metrics like replica access after creation indicate how well C3PO selected popular datasets ↔ no gauge for impact on efficient Grid processing
- Main problem: measurement of metrics related to usage of C3PO selected datasets vs. other datasets doesn't really allow for a one to one comparison

  ⇒ Decided to run C3PO in an **A/B testing mode** for a period of time

# C3PO impact analysis - A/B testing



- Direct comparison of C3PO decisions being applied vs. not being applied
  - After positive C3PO decision: coin flip (based on dataset name)
  - Decisions split into Sample A (replica is created) and Sample B (replica is not created)
  - Test period ~1.5 months
- For datasets that fall into Sample A or Sample B, compare metrics that are affected by (temporary) inaccessibility of input data or high workload on sites
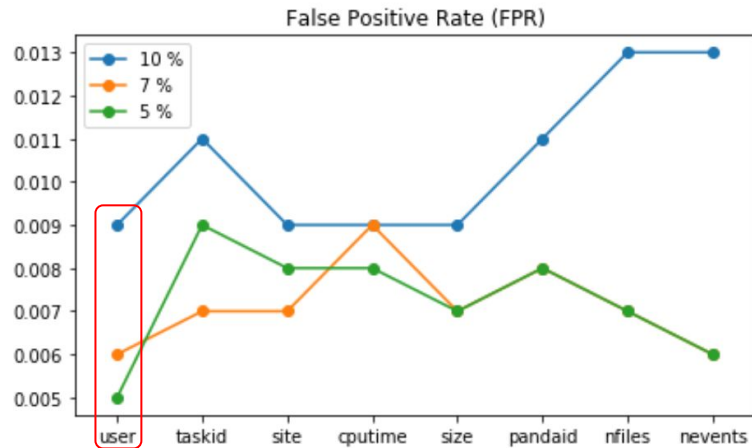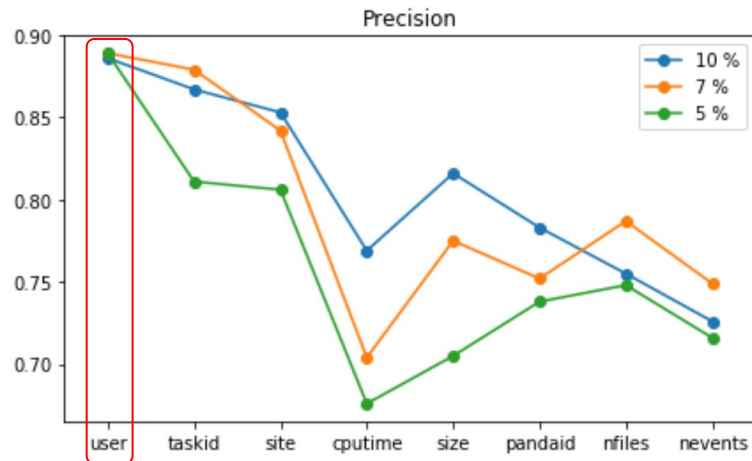
# C3PO impact analysis - Task TTC

TTC ≥ X hours

Tasks with dataset as input after C3PO decision was made
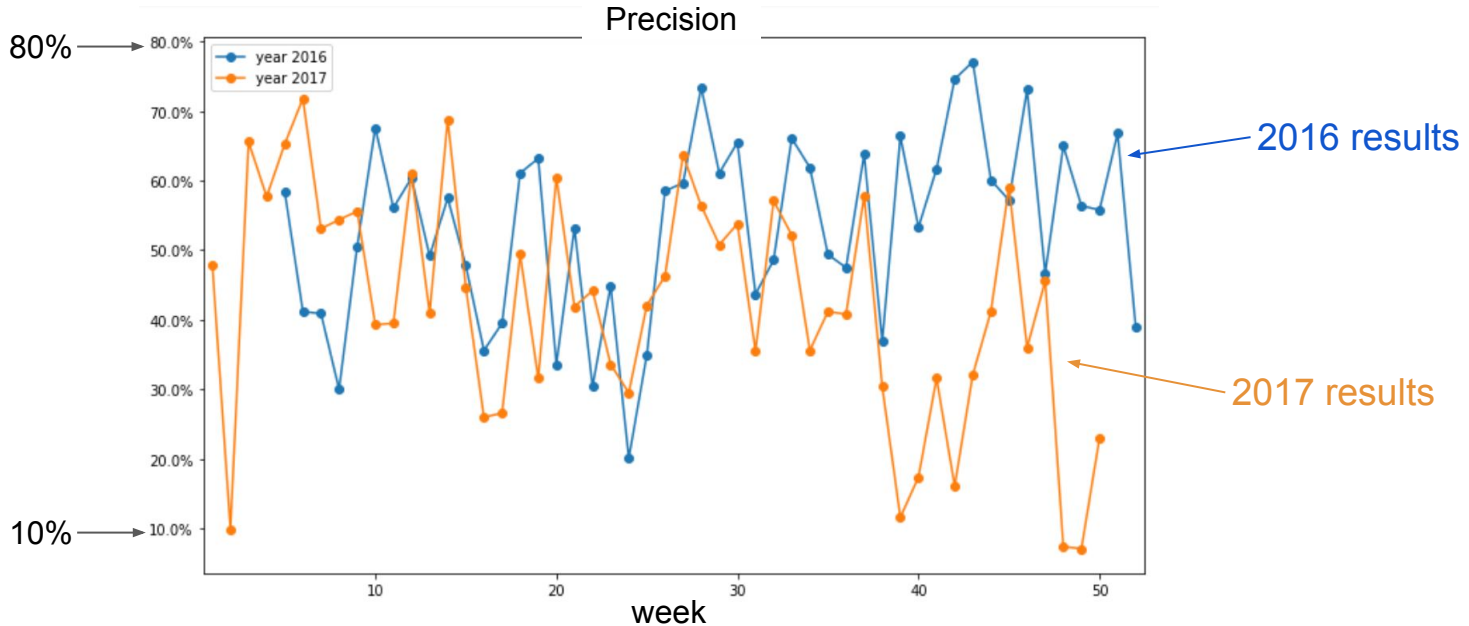


created / not created

- TTC (time to completion): time from point of creation to completion of the task
- Difference between **created** and **not created** replicas starts to occur in the tails of the TTC distribution
  - Statistically limited
  - Pronounced difference only in long tails
- Noticeable effect, but concerns only a small fraction of tasks

# Popularity prediction with machine learning (1/2)

- Bachelors project by Matteo Magoni
  - Machine learning (Adaboost decision tree) for dataset popularity prediction
  - Evaluation of historic Grid jobs meta-data
- Variables used for training
  - dsid: 6-digit dataset id
  - ptag: version tag of physics data
  - scope: data period / MC simulation campaign string
  - type: data format string
- <u>Popularity definition:</u> for a given parameter distribution, dataset falls into the tailing 10%/7%/5% (threshold cut)
- Selection of **popularity parameter**
  - evaluate data of June 2016 for several parameters (and the three threshold cuts)
  - choose **number of users with 7% cut** → highest precision (fraction of datasets predicted as popular that actually are popular)

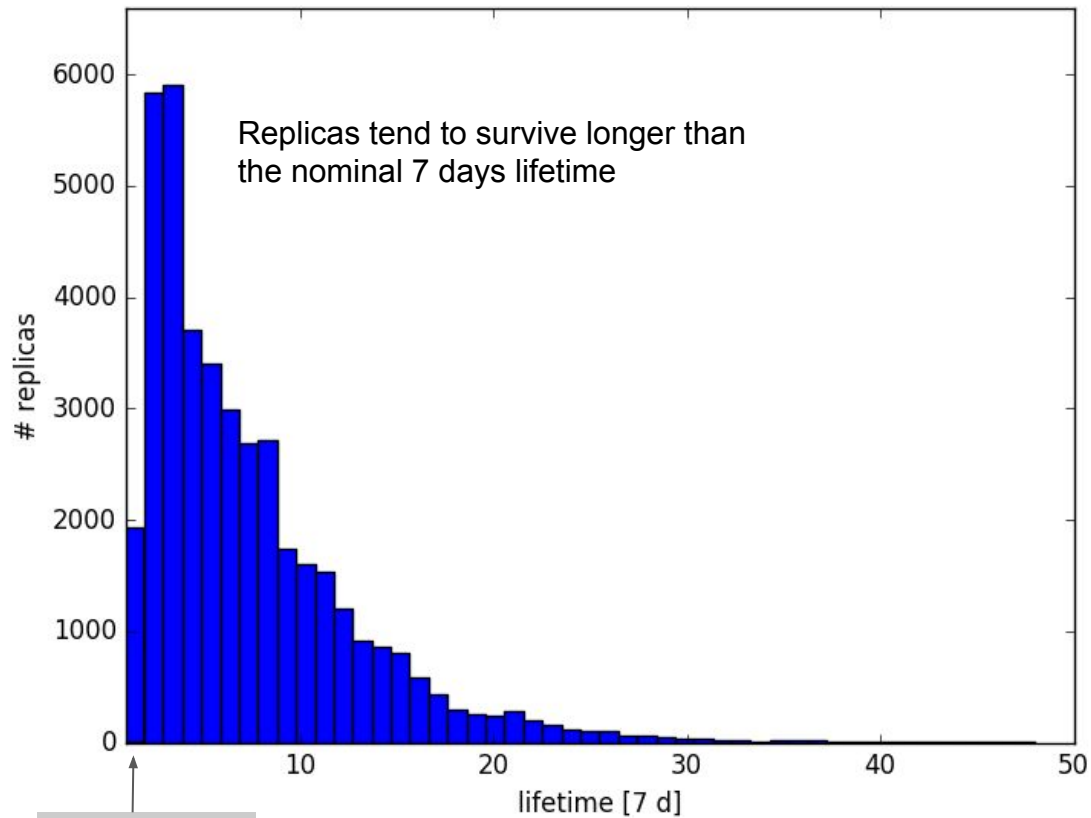# Popularity prediction with machine learning (2/2)



Precision

- Full evaluation of 2016/2017 data
  - predict popularity of given week by training on data of previous 4 weeks
  - additional training input: popularity of previous 3 weeks
- In general high resulting precision, but quite large fluctuations → lack of sufficient amount of training data?
- Several options to explore in the future (increase training window, change hyperparameters, etc.)

# Summary

- Dynamic data placement agent C3PO
  - Developed and operational during Run-2 phase of ATLAS
  - Usage efficiency of resulting newly created replicas >60% (depending on parameters like target replication site, data type/format, etc.)
  - Altering initial C3PO decision criteria affects resulting usage probability
  - C3PO replicas tend to survive longer than nominal 7 days lifetime ↔ continuously accessed → efficient use of available disk space
- C3PO impact analysis with A/B testing
  - Metrics like task TTC indicate that replicas created by C3PO have some impact on Grid processing of their associated datasets
    - Only small effects (on very limited number of Grid tasks)
    - In general difficult to unambiguously attribute observed differences to C3PO replicas
- Popularity prediction with machine learning algorithms
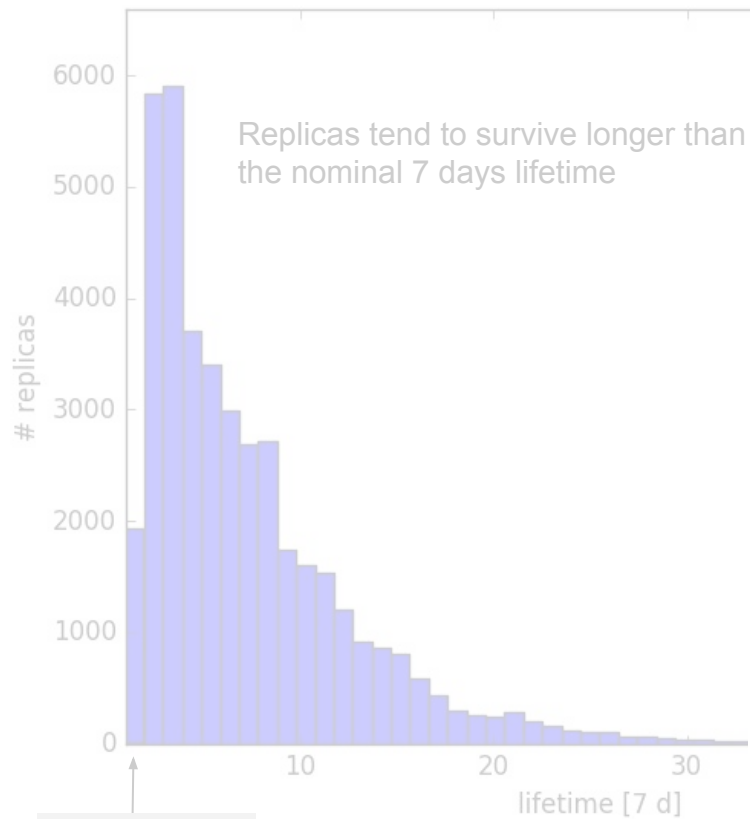  - Promising first results → multitude of options/methods to explore for improvements
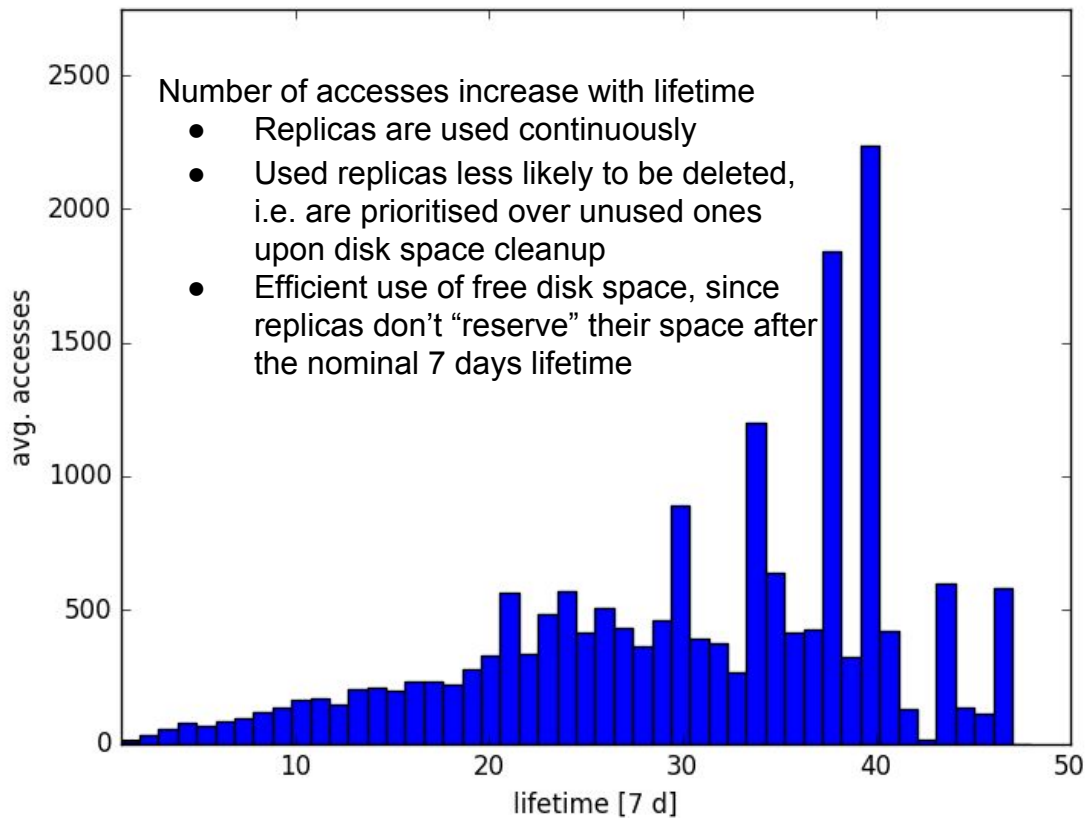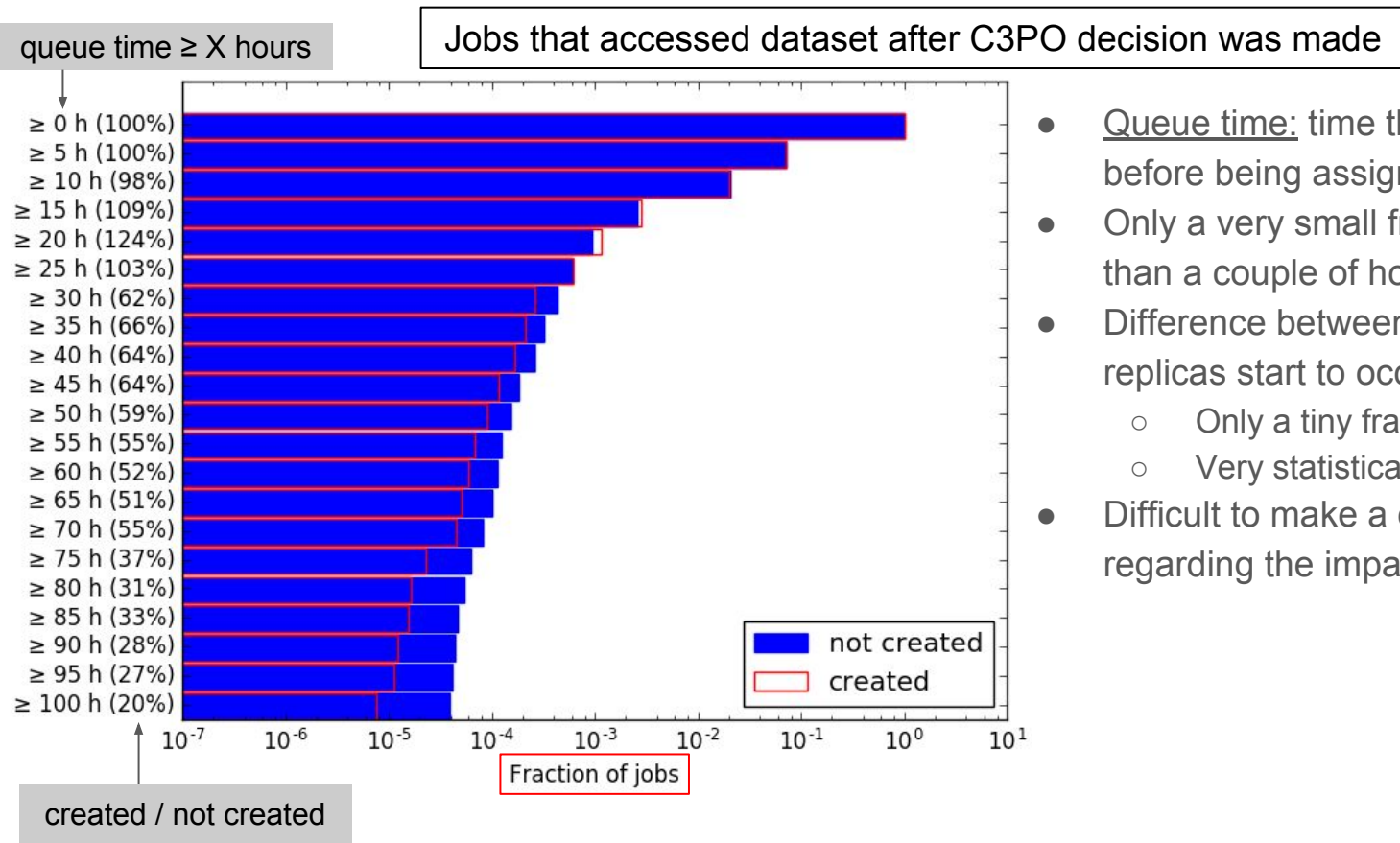
# BACKUP

# Replica lifetimes



Replicas tend to survive longer than the nominal 7 days lifetime

First 7 days

# Replica lifetimes - average number of file accesses



Replicas tend to survive longer than the nominal 7 days lifetime

# replicas

lifetime [7 d]

First 7 days

Number of accesses increase with lifetime
- Replicas are used continuously
- Used replicas less likely to be deleted, i.e. are prioritised over unused ones upon disk space cleanup
- Efficient use of free disk space, since replicas don't "reserve" their space after the nominal 7 days lifetime

avg. accesses

lifetime [7 d]

# C3PO impact analysis - job queue times



queue time ≥ X hours

Jobs that accessed dataset after C3PO decision was made

created / not created

- Queue time: time that job waits in site queue before being assigned to worker node
- Only a very small fraction of jobs wait longer than a couple of hours in queue
- Difference between **created** and **not created** replicas start to occur in the long tails
  - Only a tiny fraction of jobs
  - Very statistically limited!
- Difficult to make a decisive conclusion regarding the impact on job queue times