# Event Driven Processing and Data Management

Paul Millar
Joint ESCAPE WP2/WP5 workshop
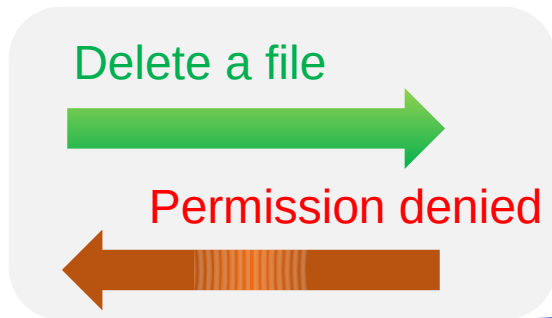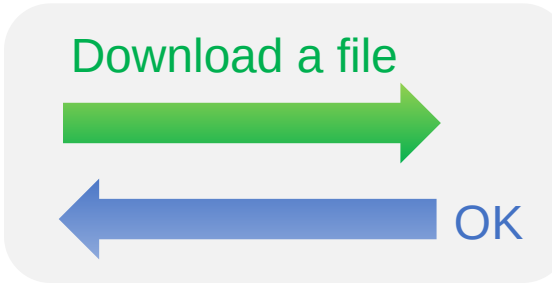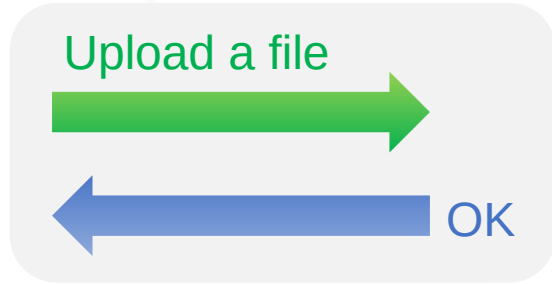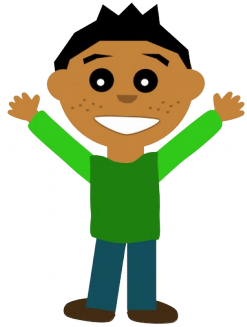
(with material donated by Marica Antonacci, Patrick Fuhrmann and Michael Schuh)
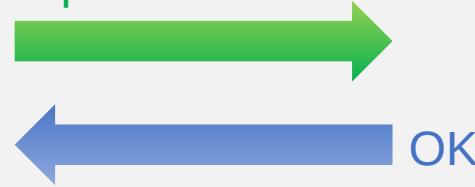
# How storage is used currently

# This may lead to problems...

**Rucio/LFC/...**

**Storage Node**

Upload a file

OK

Register file
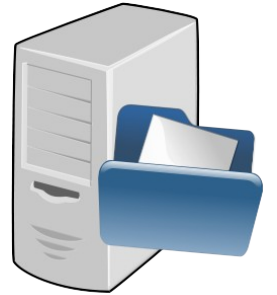
OK

Delete a file

OK

Unregister file

OK

Leads to
- Dark Data
- Dangling References

# This may cause more problems

Stage file A, file B, ..... , file X

Requests Queued

Are files on DISK ?

Yes, no, Yes ... No

Are files on DISK ?

Yes, Yes, Yes... No

\* \* \*

Are files on DISK ?

Yes, Yes, ... ERROR

# New way of interacting: storage events

Storage Events

File uploaded

File moved to tape

Storage Media Transitions

Tape

SSD

DISK

File staged to disk

Admin intervention

(e.g., data lost)

# dCache implementation

# dCache Storage Events: Kafka and SSE
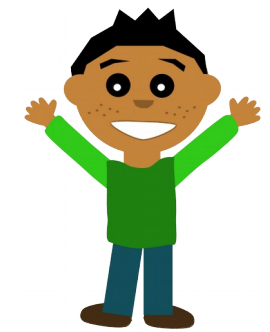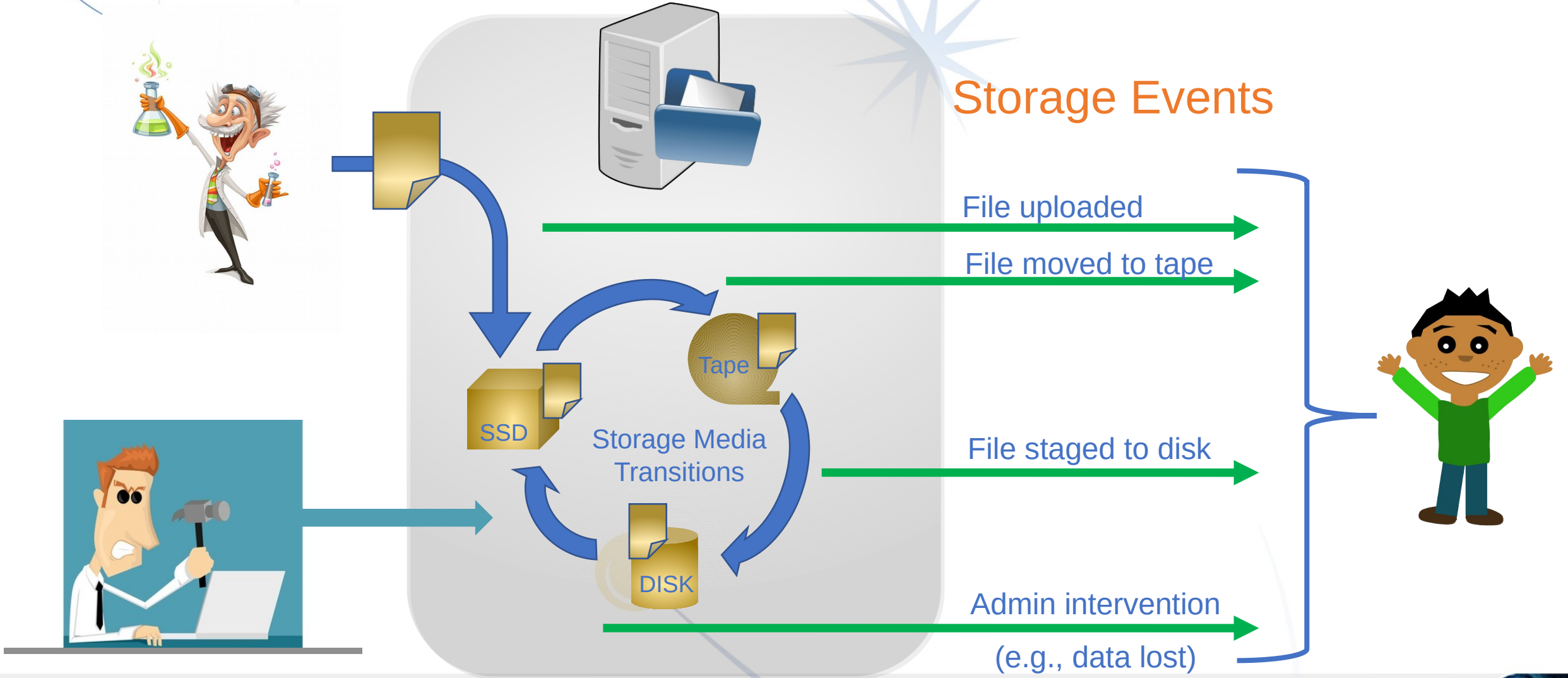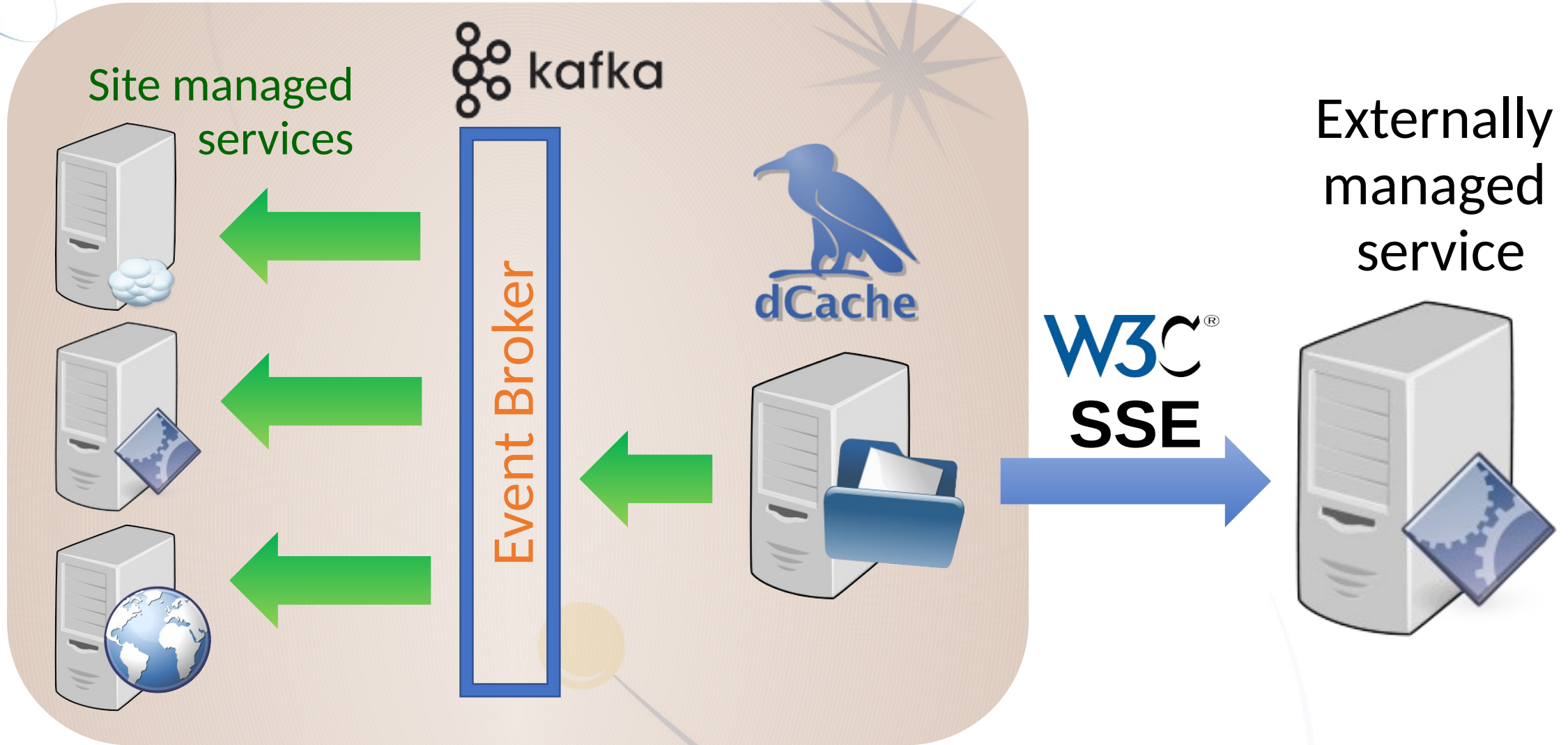
# Cheat sheet: Kafka vs SSE

| | kafka | W3C® SSE |
|---|---|---|
| **Availability since...** | dCache v4.1 | dCache v5.0 |
| **Standard …** | Software package | Protocol |
| **What events does it see?** | dCache billing events | inotify |
| **Main benefit** | Easy integration | Built-in security |
| **"Catch-up" storage** | Memory & disk | Memory-only (currently) |
| **Target audience** | Site-level integration | Events for users |

# Use-cases and demonstrators

# INDIGO-Orchestrator
# & automated data processing

# INDIGO Orchestrator (SSE)

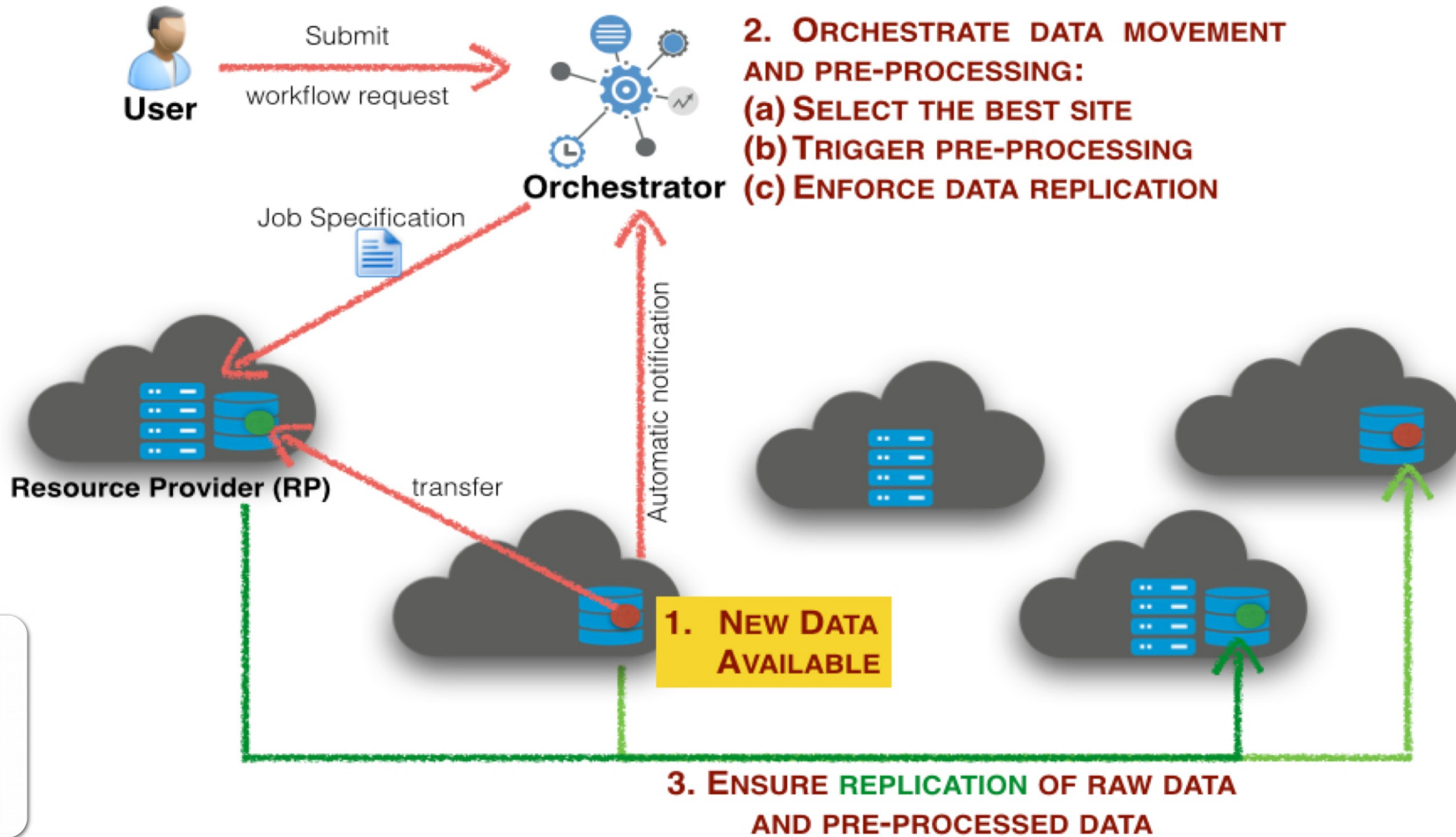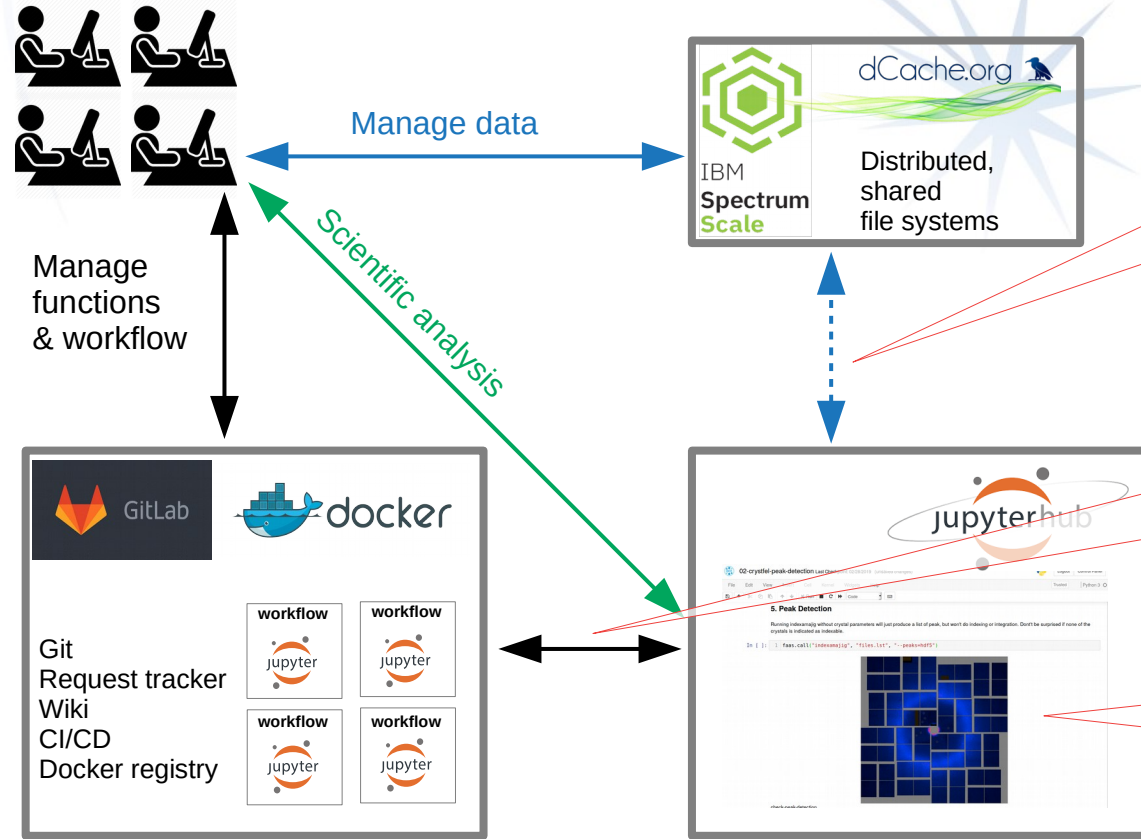# Use-cases and demonstrators

# EU-XFEL
# Analysis and automation pilot platform pilot

# Jupyter Notebooks on HPC cluster



- Data on shared file systems mounted on HPC host (NFS)
- uid, gid mapping

Sync git and Jupyter
- Mybinder
- Nbgitpuller
- git

Reproducibility
- Functions
- Results

Manage data

Manage functions & workflow

Scientific analysis

Distributed, shared file systems

Git
Request tracker
Wiki
CI/CD
Docker registry

workflow
workflow
workflow
workflow

→ user action    ----▶ program action    ● Data    ● Code    ● Analysis

Contact: eosc-pan-info@desy.de
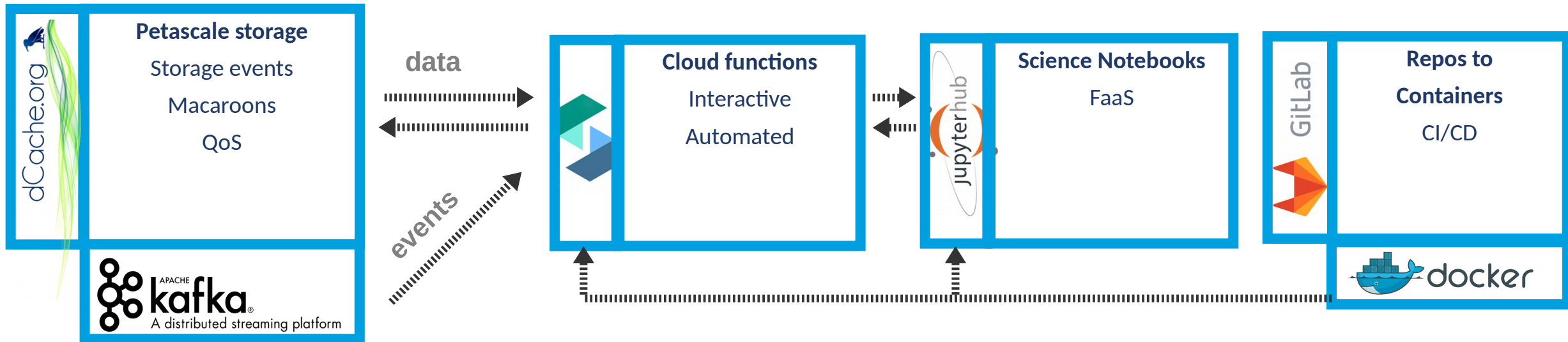Icons: flaticon.com (freepik/prettycon)

# Analysis and automation platform

# Analysis and automation platform
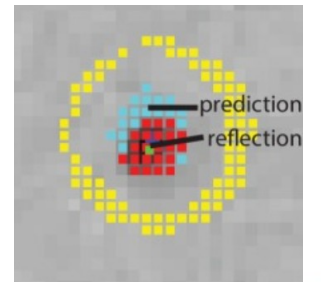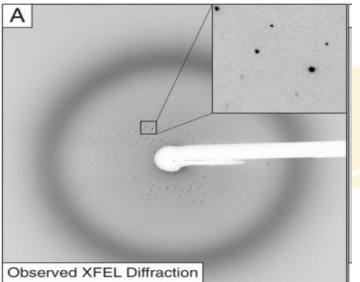
Single namespace
in multi-clouds.

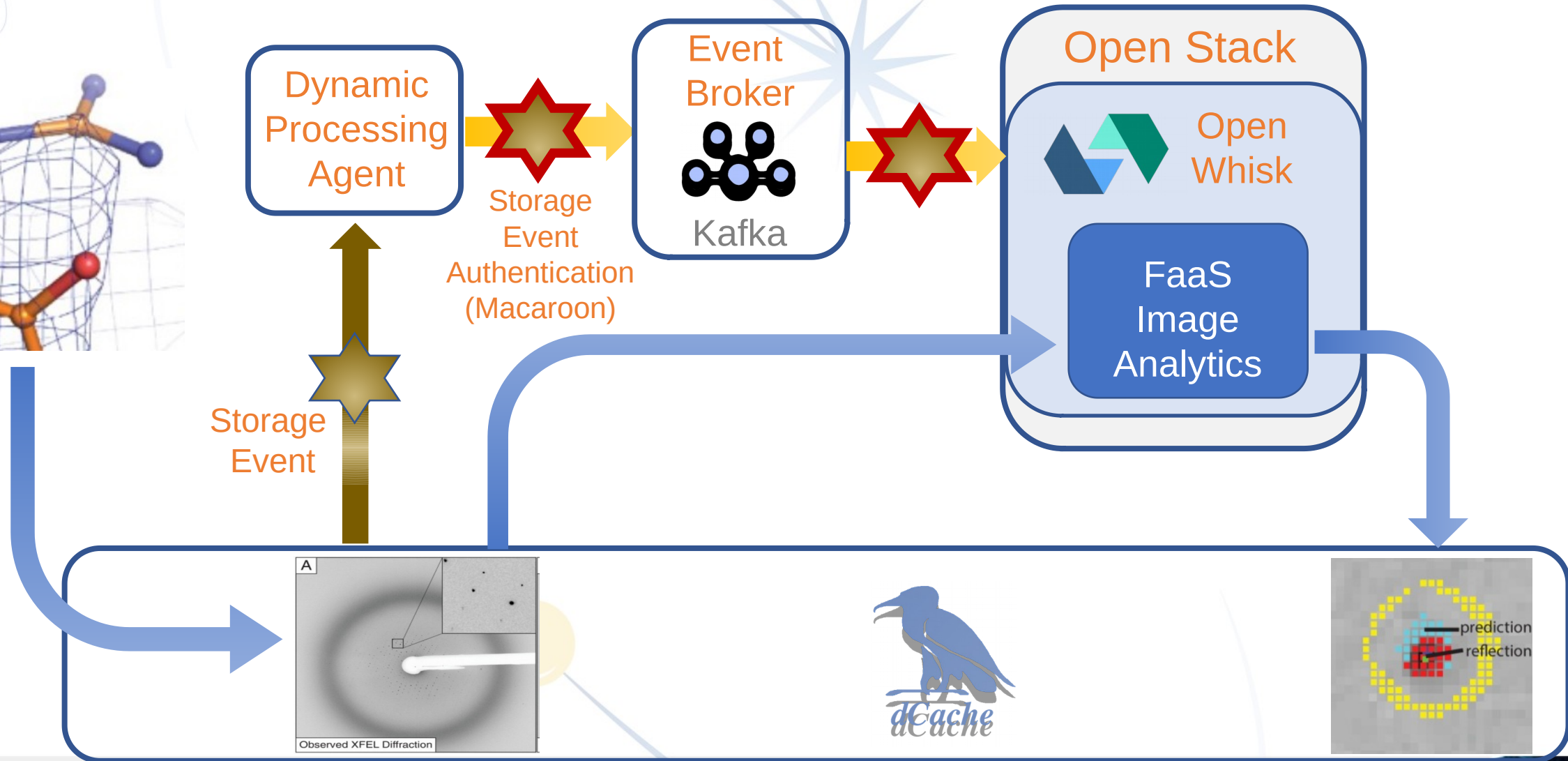Function-as-a-Service
in Science Notebooks
and in automation.

Jupyter Notebooks
in user-defined
environments.

Just push code
it builds, goes live
and scales.

**Petascale storage**

Storage events

Macaroons

QoS

data

events

**Cloud functions**

Interactive

Automated

**Science Notebooks**

FaaS

**Repos to Containers**

CI/CD

# System/integration view

# EISCAT 3D
# Automated replication pilot

# EISCAT_3D: automated replication

**dCache**

**Register for events**

**W3C SSE**

**Rucio Panoptes Agent**

**RUCIO**

# EISCAT_3D: automated replication

Upload a file

dCache

Rucio Panoptes Agent

RUCIO

01/07/2019

# EISCAT_3D: automated replication



**dCache**

**New file event**

**W3C® SSE**

**Rucio Panoptes Agent**

**RUCIO**

# EISCAT_3D: automated replication



dCache

**Rucio Panoptes Agent**

**New file event**

**Rucio internal bus**

RUCIO

# EISCAT_3D: automated replication

**Rucio Panoptes Agent**

dCache

Evaluate rules. Choose second replica location.

RUCIO

# EISCAT_3D: automated replication



dCache

**Rucio Panoptes Agent**

RUCIO

Evaluate rules. Choose second replica location.

# EISCAT_3D: automated replication

dCache

**Rucio Panoptes Agent**

RUCIO

FTS

**Initiate a third-party copy**

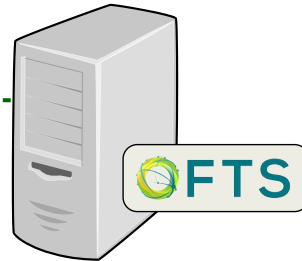# EISCAT_3D: automated replication



dCache

**Rucio Panoptes Agent**

RUCIO

**Third-party copy**

FTS

# EISCAT_3D: automated replication

01/07/2019

# Use-cases and demonstrators

# Fermilab & WLCG
# Increased tape staging efficiency

# Increased tape staging efficiency

- One of the closely correlated metric for tape inefficiency is **remount count**: loading same tape two (or more) times in close succession.

  Load tape once. Read all relevant data. Move on.

- To reduce likelihood, the tape system should be given (ideally) **all pending requests**.

  This allows the tape system to reorder requests.

- **Problem**: polling overhead from repeatedly checking status.

- **Solution**: storage events to discover when files have been staged

  - Allow clients to request all files (almost without limit).

  - Low latency: jobs/transfers may start as soon as file is available.

# Increased tape staging efficiency

Stage file A, file B, ..... , file X

Requests Queued

File 1 arrived on DISK

File 2 arrived on DISK

File 3 arrived on DISK

*  *  *

File N arrived on DISK

# Handling data-loss:
# Automated data-replica recovery

# Automatic data-replica recovery



**dCache**

**Register for events**

**W3C® SSE**

**Rucio Panoptes Agent**

**RUCIO**

# Automatic data-replica recovery



**dCache**

**Rucio Panoptes Agent**

**RUCIO**

# Automatic data-replica recovery

**Rucio Panoptes Agent**

dCache

# Automatic data-replica recovery



**dCache**

**File lost event**

**W3C SSE**

**Rucio Panoptes Agent**

**RUCIO**

# Automatic data-replica recovery

**dCache**

**Rucio Panoptes Agent**

**Replica lost**

**Rucio internal bus**

# Automatic data-replica recovery

dCache

**Rucio Panoptes Agent**

Find "best" storage with a replica of this file

RUCIO

# Automatic data-replica recovery



dCache

Rucio
Panoptes
Agent

RUCIO

FTS

**Initiate a
third-party copy**

# Automatic data-replica recovery



**dCache**

**Rucio Panoptes Agent**

**RUCIO**

**Third-party copy**

**FTS**

# Automatic data-replica recovery

# Thanks for listening!

# Bonus material

# New solutions to old problems



Upload a file
OK

Delete a file
OK

Storage Node

Register file
OK

Unregister file
OK

Rucio/LFC/...

# New solutions to old problems

Stage file A, file B, ..... , file X

Requests Queued

File 1 arrived on DISK

File 2 arrived on DISK

File 3 arrived on DISK

* * *

File N arrived on DISK

43

# Comparison: events in industry...

## Amazon Lambda

### Image Thumbnail Creation



*Photograph is taken* ... *Lambda is triggered*

S3

Photo is uploaded into S3 bucket

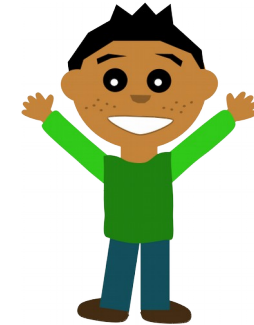Lambda runs image resizing code to generate Web, mobile and tablet sizes

## Google Cloud Platform

# Comparison: events in Open-Source

**STORM**

Apache Storm is a distributed stream processing computation framework written predominantly in the Clojure programming language.

**Spark**

**OpenWhisk**

**APACHE kafka®**
A distributed streaming platform

**samza**

Samza allows you to build stateful applications that process data in real-time from multiple sources including Apache Kafka.

**APACHE nifi**

Apache NiFi is a software project from the Apache Software Foundation designed to automate the flow of data between software systems.
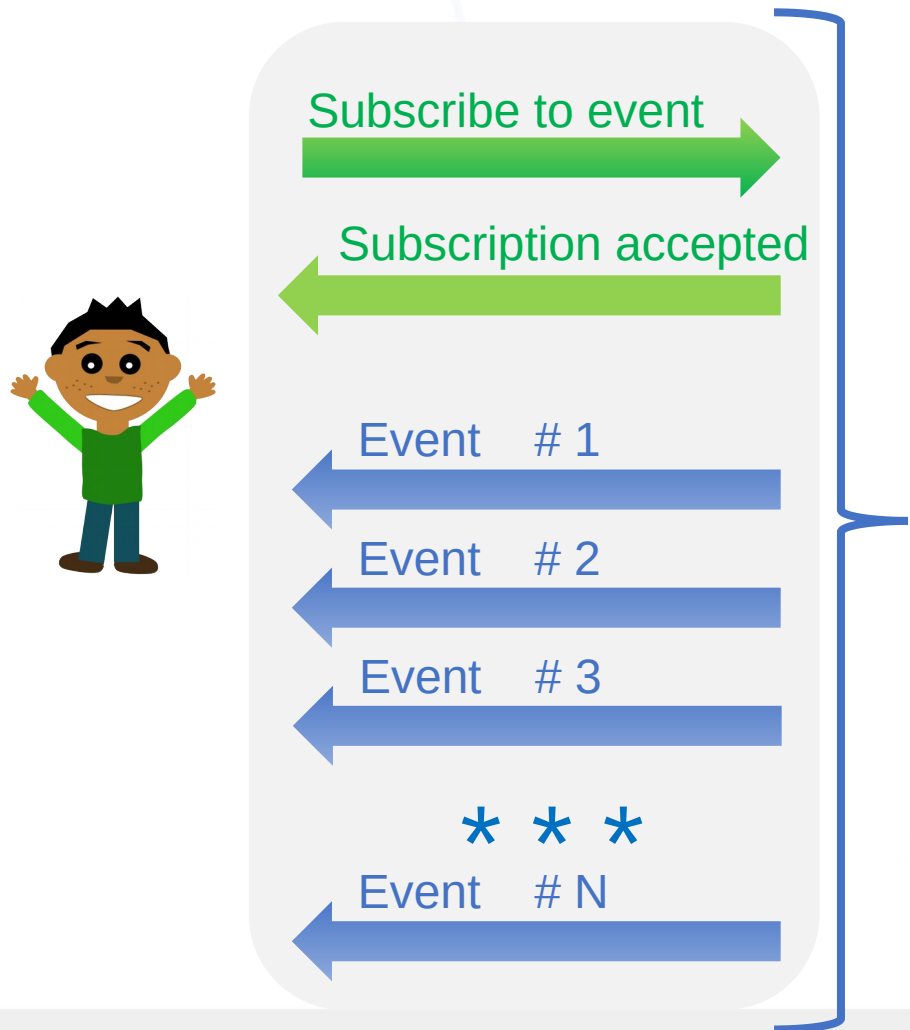
**Kubeless**

Kubeless is a Kubernetes-native serverless framework that lets you deploy small bits of code (functions) without having to worry about the underlying infrastructure.
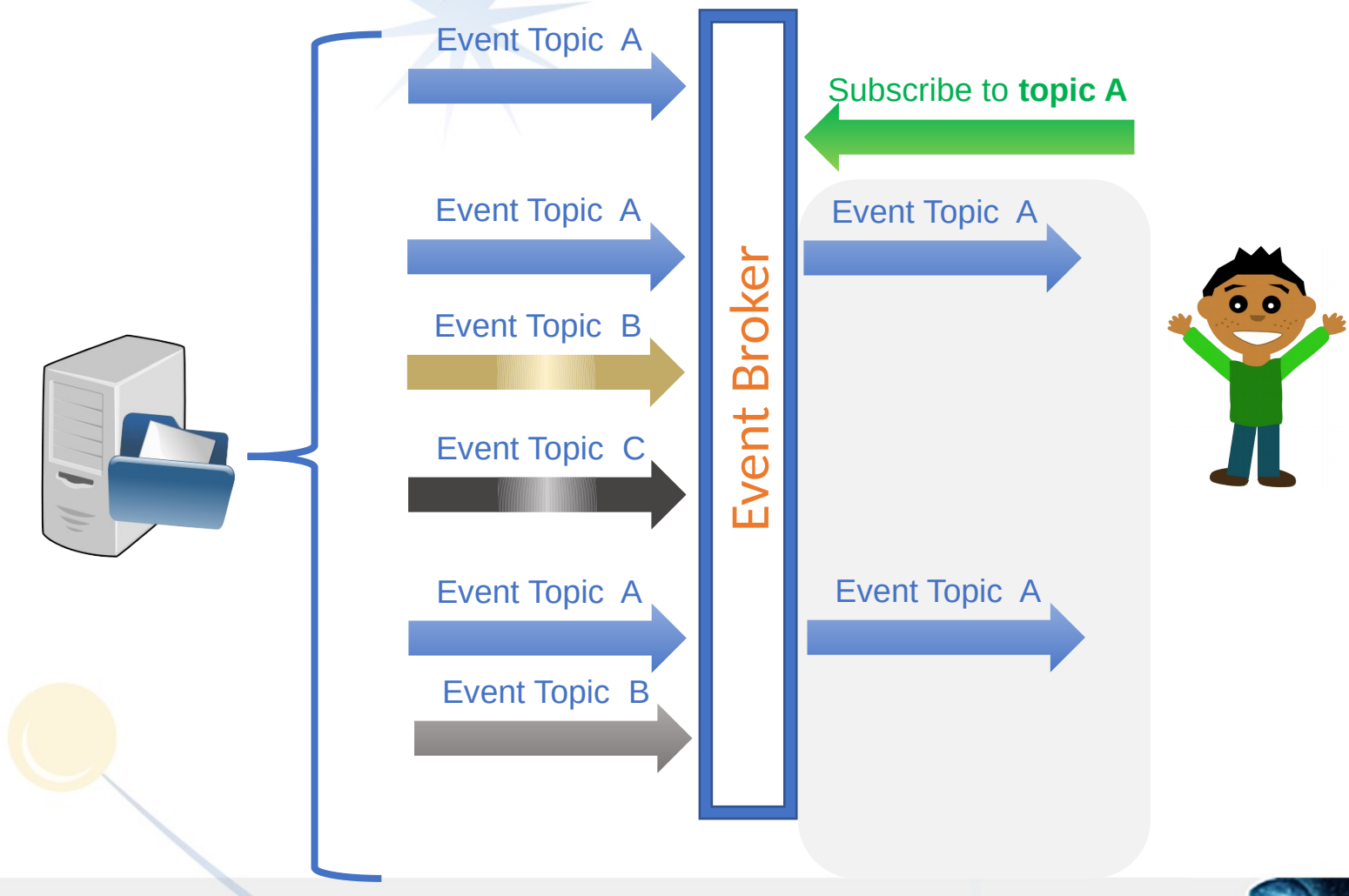
# New way of interacting: storage events

## Direct event delivery

Subscribe to event

Subscription accepted

Event # 1

Event # 2

Event # 3

* * *

Event # N

## Brokered event delivery

Event Topic A

Event Topic A

Event Topic B

Event Topic C

Event Topic A

Event Topic B

Event Broker

Subscribe to **topic A**

Event Topic A

Event Topic A

# Full XDC Layout using Storage Events