# USABILITY

# +COLLABORATION

# Agenda

- Border area
- [DEMO] Recycling YAML files
- Cars and washing machines
- [DEMO] Collaboration and Usability

# Márton Elek

- Apache Hadoop/Ratis PMC
  - elek@apache.org
- Cloudera
  - Principal software engineer
- twitter.com/@anzix
- Ozone:
  - Ozone Explained (Youtube)
- Kubernetes + Apache Bigdata:
  - **github.com/elek**/flekszible
  - flokkr.github.io

# Márton Elek

- Apache Hadoop/Ratis PMC
  - elek@apache.org
- Cloudera
  - Principal software engineer
- twitter.com/@anzix
- Ozone:
  - Ozone Explained (Youtube)
- Kubernetes + Apache Bigdata:
  - **github.com/elek**/flekszible
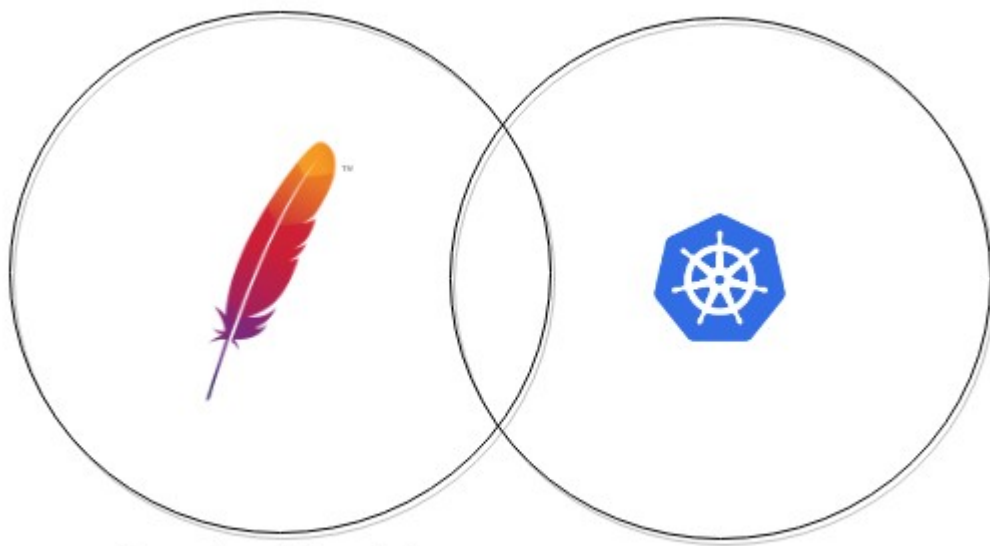  - flokkr.github.io

# How to manage cluster?

# Borders

# Nobody's land
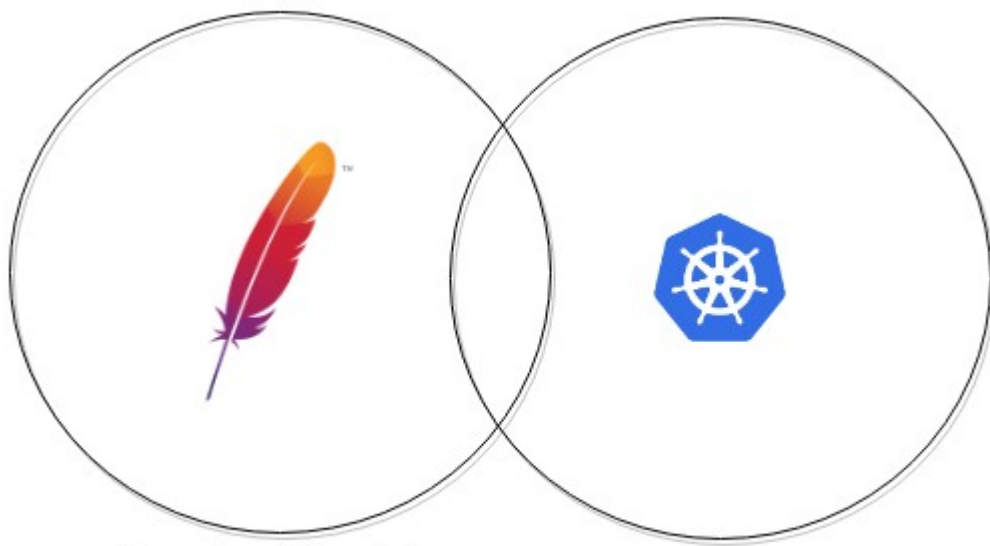
Apache Big-data       Cloud native

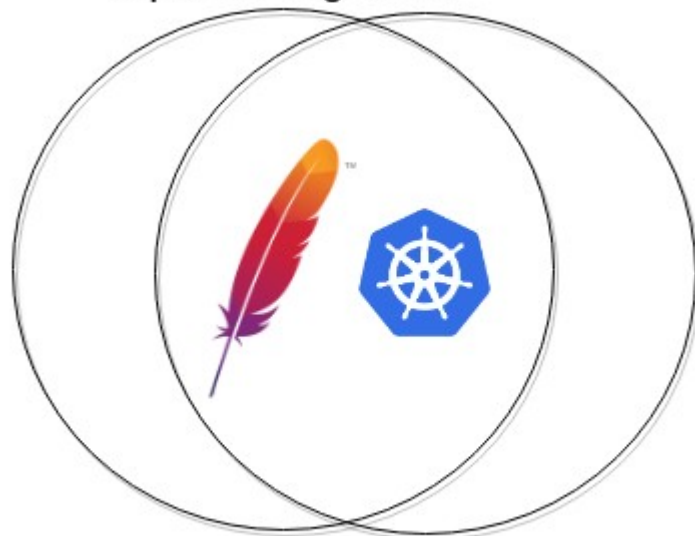Apache Big-data

Cloud native

Apache Big-data                    Cloud native
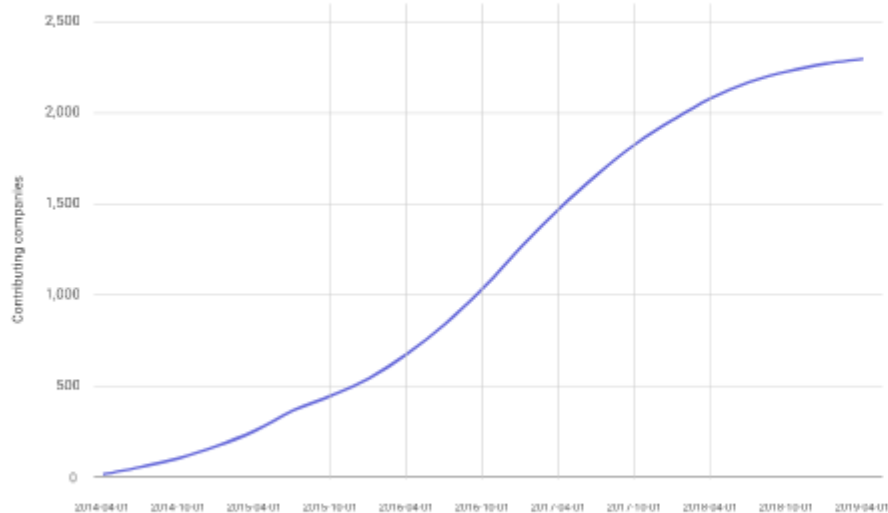
Apache Big-data



Cloud native

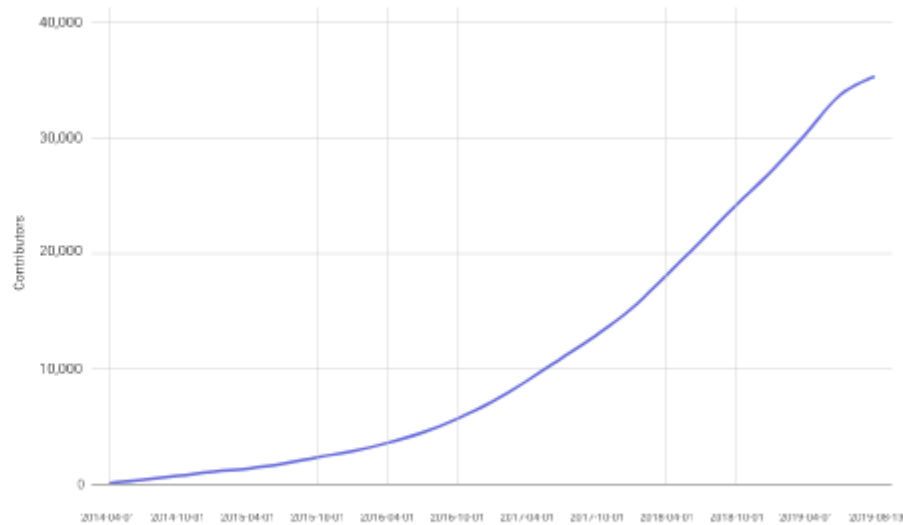# Cloud-native

# Why is it so popular?

## Growth: Contributing Companies and Contributors



*Cumulative growth of number of contributing companies*



*Cumulative growth of contributors*
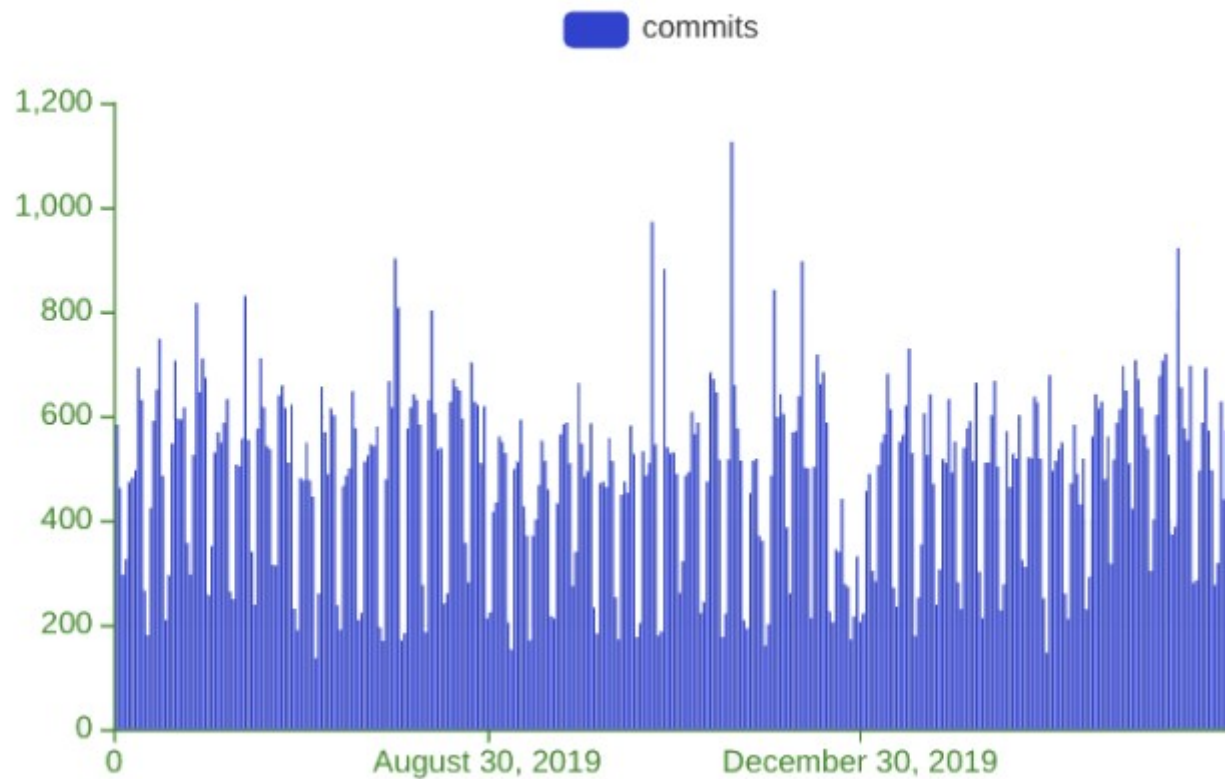
# Why is it so popular?

- Good low level design
- Flexible and generic api
  - apiVersion/kind → CRUD
- YAML files everywhere
- → Very popular

# COLLABORATION

# Apache Big-data

# Apache ...

**Commit History**

# Apache ... Big-Data

Number of resolved issues in Jira (FY21: 2019-05-01 - 2020-05-01)

(Hadoop projects are combined)

...

| project | key |
|---|---|
| FLINK | 3040 |
| HADOOP | 2819 |
| ARROW | 2562 |
| SPARK | 2419 |
| AIRFLOW | 1846 |
| INFRA | 1707 |
| BEAM | 1464 |
| CAMEL | 1324 |
| HBASE | 1305 |
| GEODE | 1060 |
| HIVE | 996 |

# Apache Big-data ecosystem

- It's an ecosystem
  - Collaboration is in the genes
- Hadoop compatible file system
  - Used by Spark, Flink, Hive, Hbase…
- Hadoop RPC/Configuration
  - Used by many projects

# COLLABORATION

# Nobody is an island

- Most of big-data projects used **together** with others
  - Kafka alone?, Spark alone?
- There are only a few meta-project
  - Apache Ambari: cluster management
  - Apache Bigtop: packaging in the same way

# COLLABORATION?

# Collaboration

- Why don't we have more collaborative, meta-projects?

- Projects are complex and hard to configure them

    - Complexity is increasing with each project

- Can K8s help in the collaboration?

# USABILITY?

# Recycling

# K8s resource management?

- Helm
  - No real collaboration, not flexible
- Kustomize
  - Transformations are not re-usable
- Any new tool?
  - Reusable components + transformations

# github.com/elek/flekszible

- Kubernetes Yaml file manager
  - Reusable k8s files
  - **Reusable transformations (!!!)**
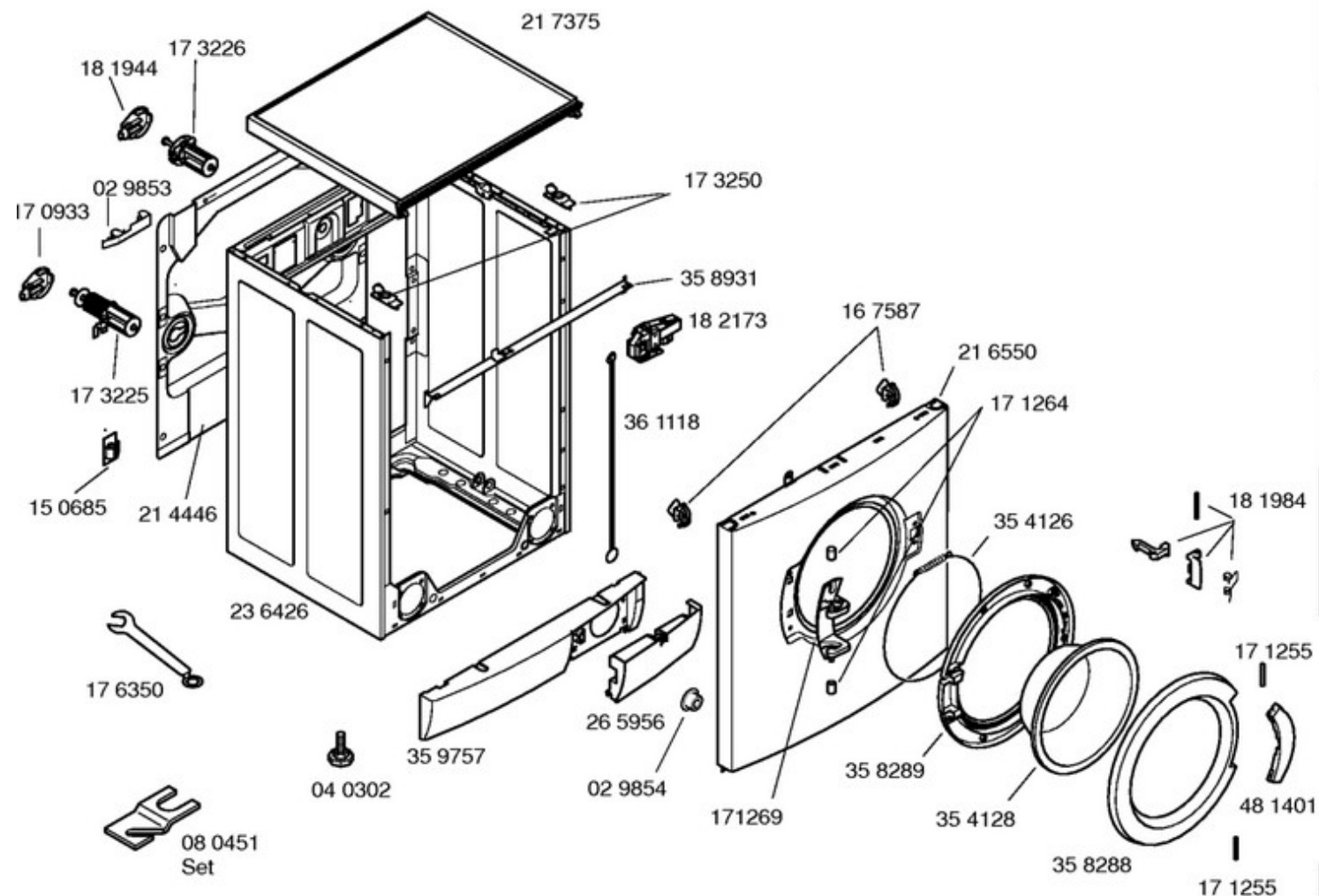- → Final cluster definition

DEMO

# flokkr.github.io

- Docker images for Apache big-data projects

- Reusable k8s files + **transformations!!**

- Did we solve all the problems:
  - COLLABORATION: DONE
  - USABILITY: ???

# Config granularity

**Dämmatten–Set**
**21 6336**

21 7375

18 1944

17 3226

17 0933

02 9853

17 3250

17 3225

35 8931

18 2173

16 7587

21 6550

171 264

36 1118

150 685

214 446

18 1984

236 426

354 126

176 350

359 757

265 956

18 1984

171 255

040 302

080 451
Set

029 854

171 269

358 289

354 128

481 401

358 288

171 255

# Usability

- The challenge is to find the right abstraction level
  - [Top] Give me a secure HDFS+Kafka
  - ???
  - [Bottom]: HDFS-SITE.XML_dfs.checksum.type: "NULL"

# Cluster management

# How to manage cluster?

# COLLABORATION
# +USABILITY

DEMO

# Summary

- Kubernetes is a good tool to improve **Collaboration** with increased **Usability**

- But we need good tool with the good abstraction

- New Apache "meta-projects" can be useful
  - → Incubator? (if interested → elek@apache.org)

# Márton Elek

- elek@apache.org
- twitter.com/anzix
- Ozone:
  - Ozone Explained (Youtube)
- Kubernetes + Apache Bigdata:
  - **github.com/elek**/flekszible
  - flokkr.github.io