# Adding intelligence to stream-processing applications

Michael McCune and William Benton

msm@redhat.com • @FOSSJunkie

willb@redhat.com • @willb

# Forecast

Introducing intelligent applications

Introducing microservice architectures

Stream processing with Apache Kafka and Apache Spark
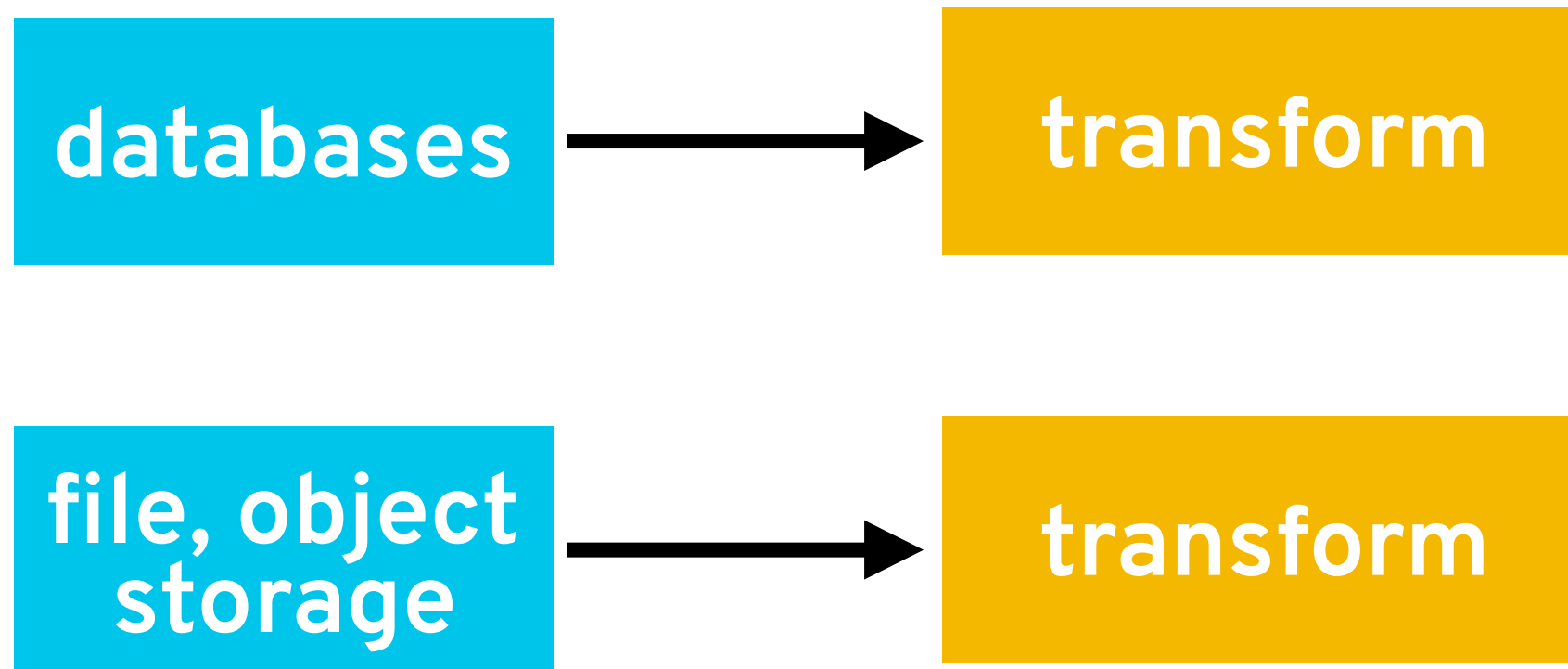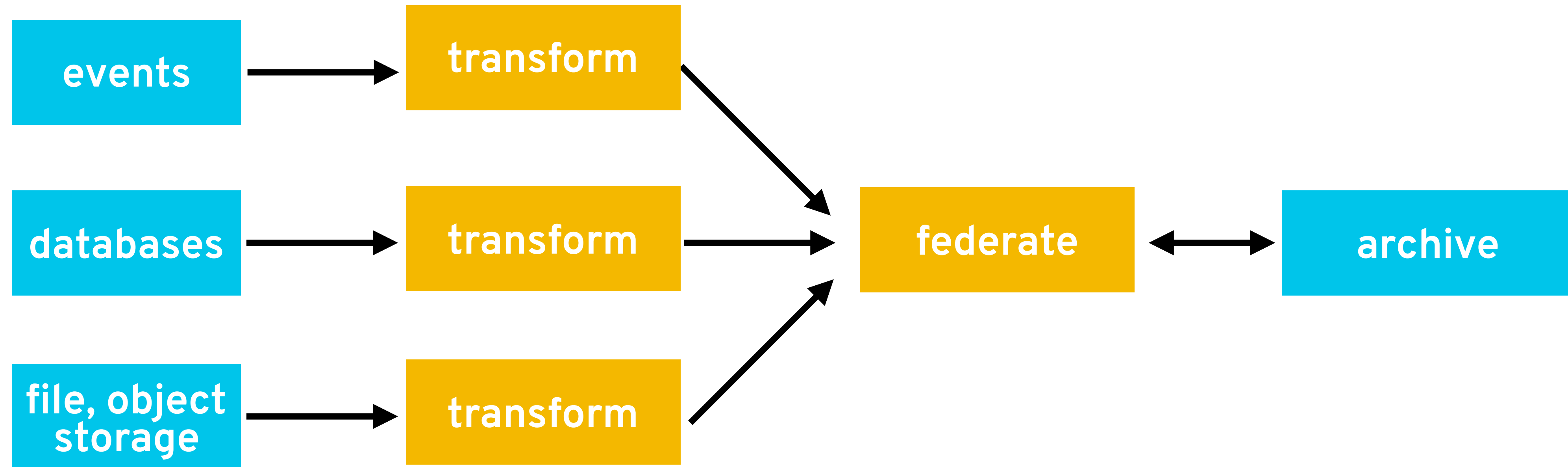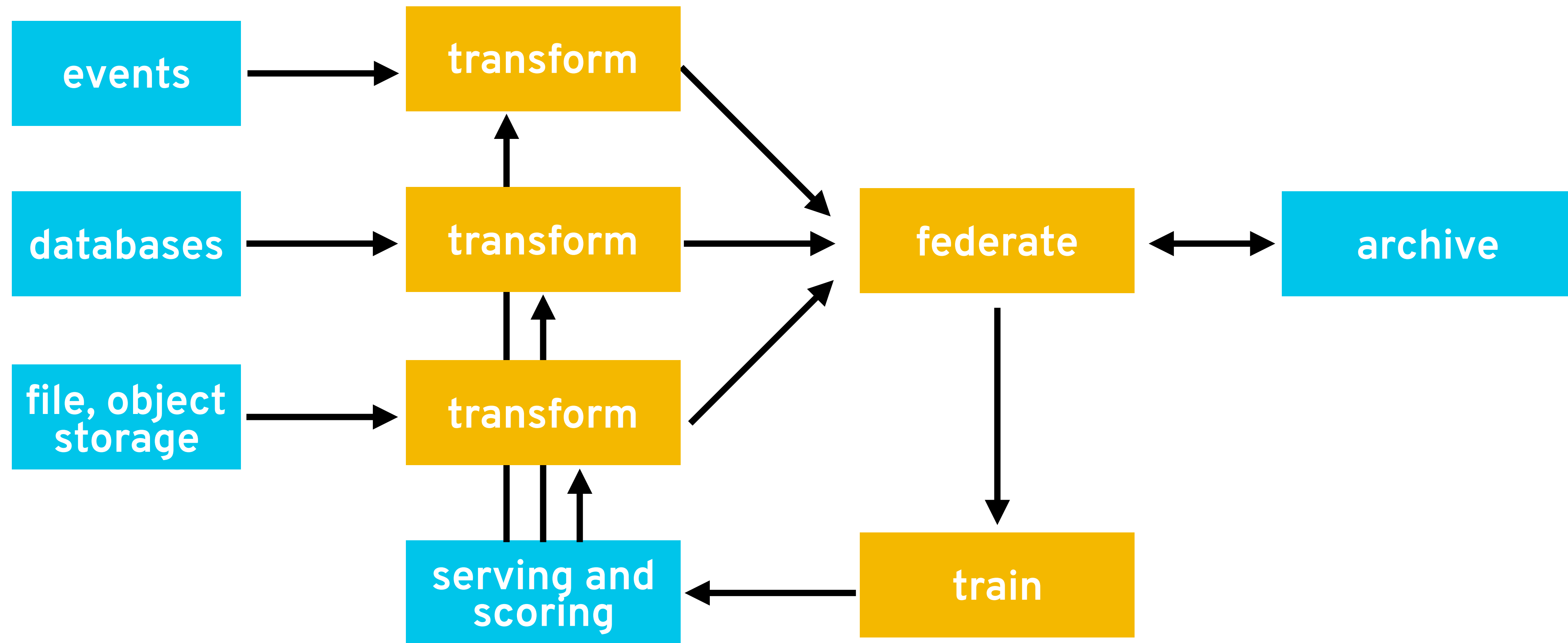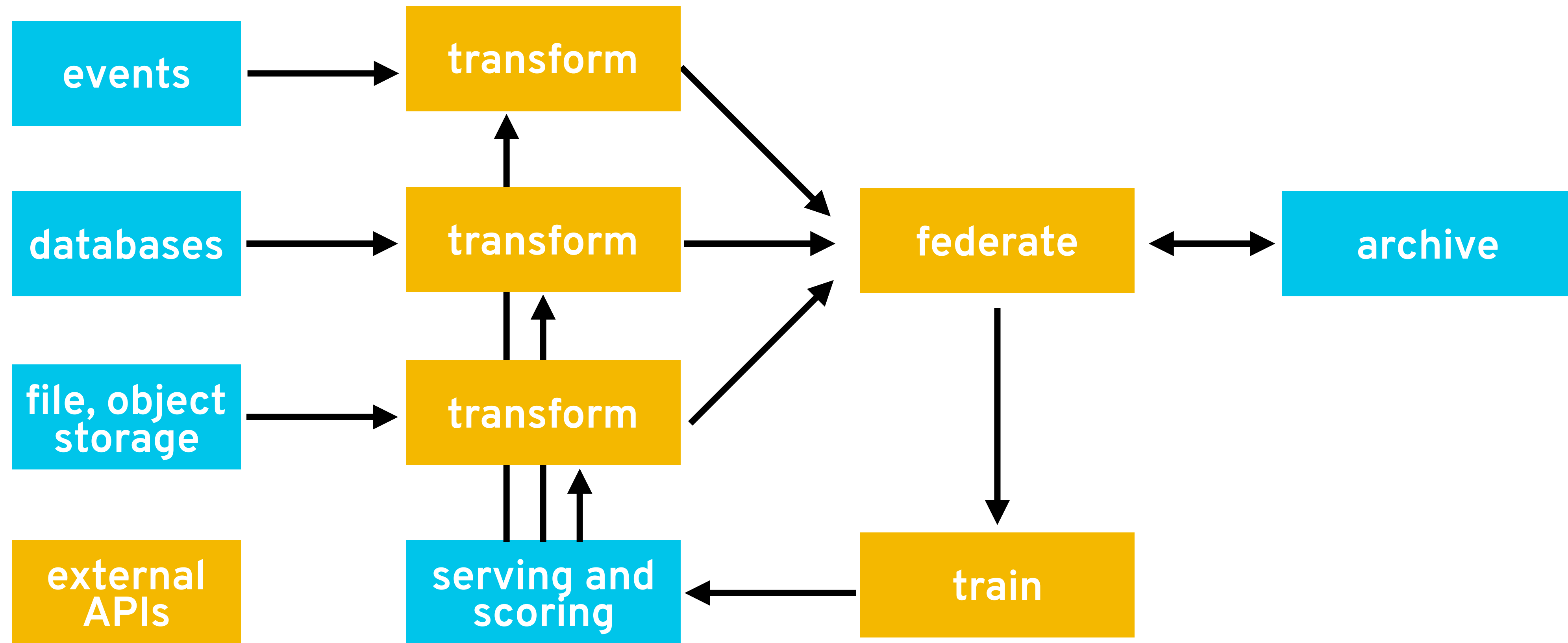
Interactivity time

Intelligent applications **collect** and **learn from** data in order to provide **improved functionality** with **longevity** and **popularity**.
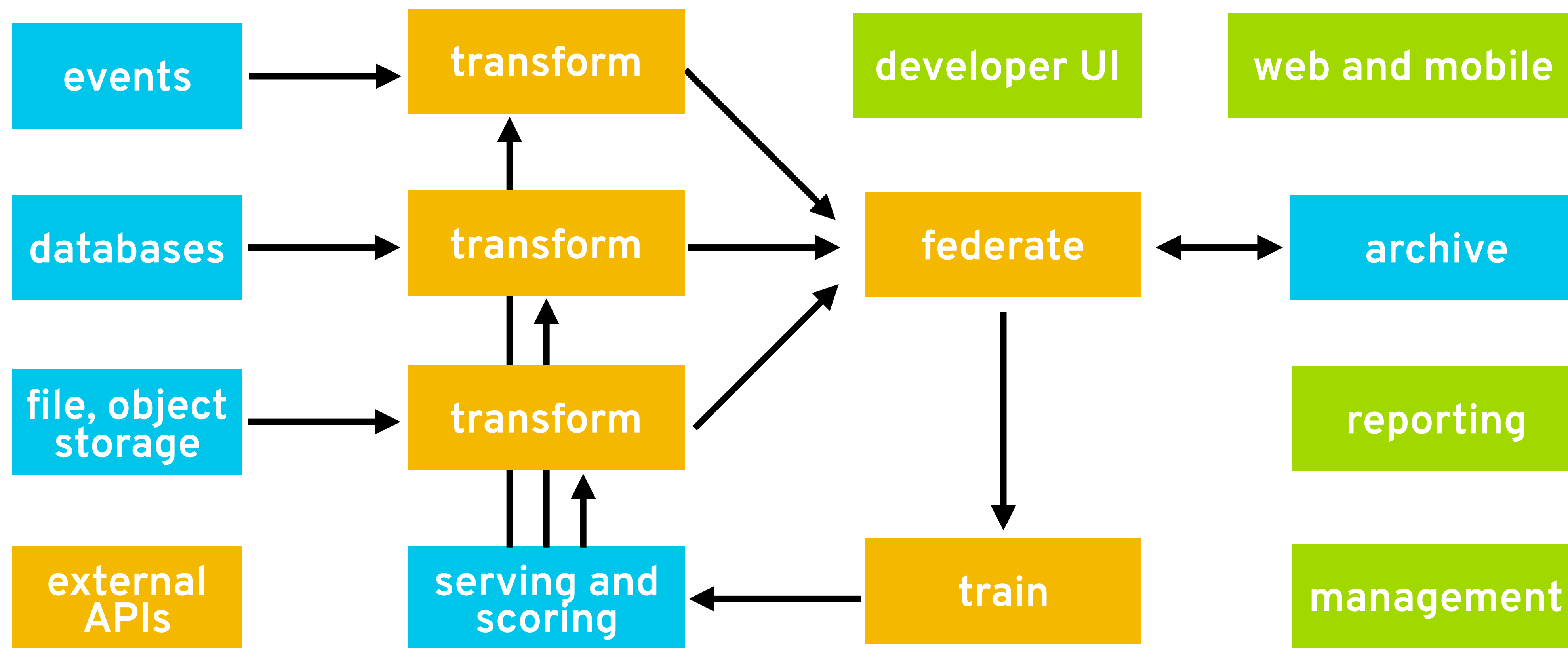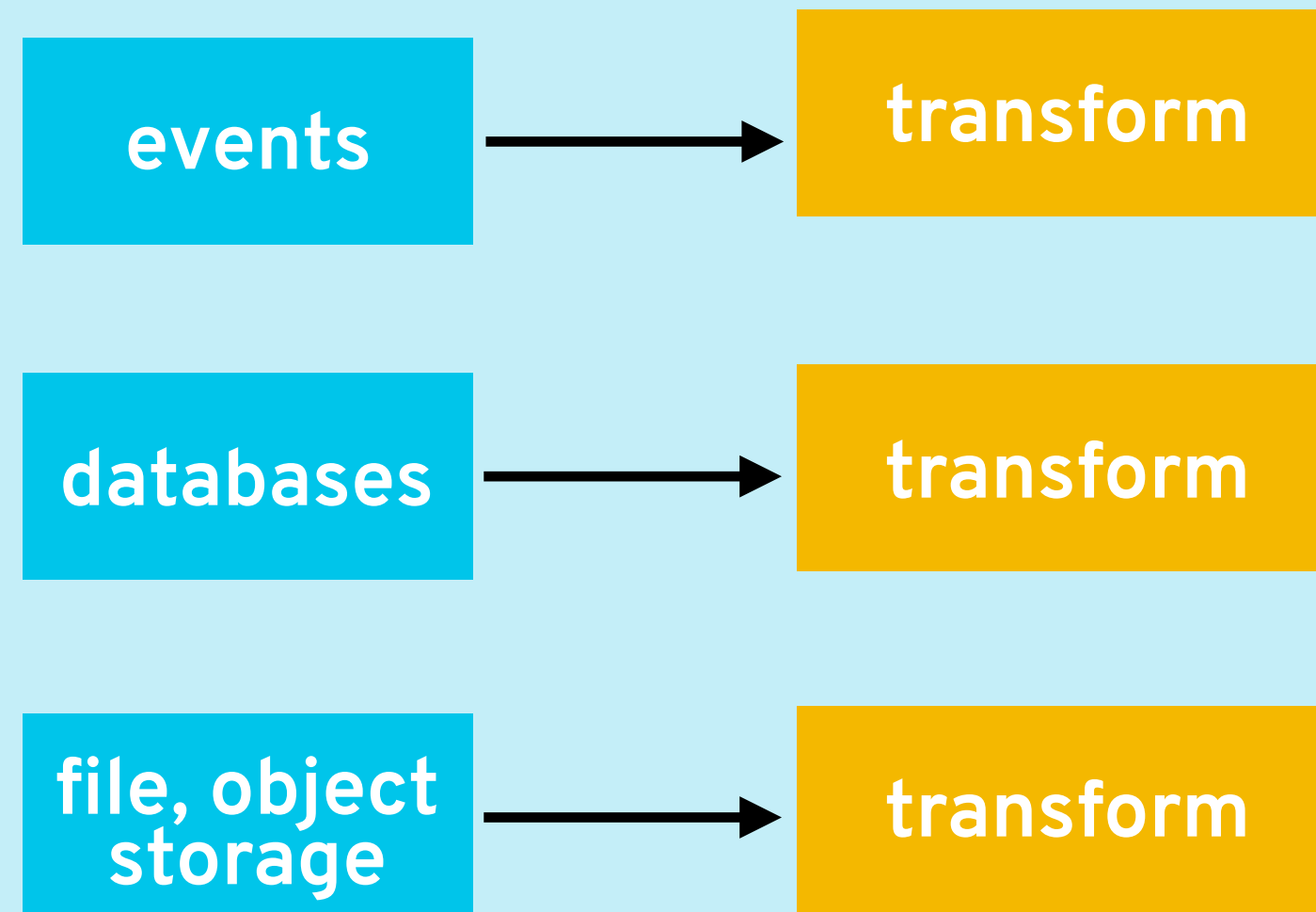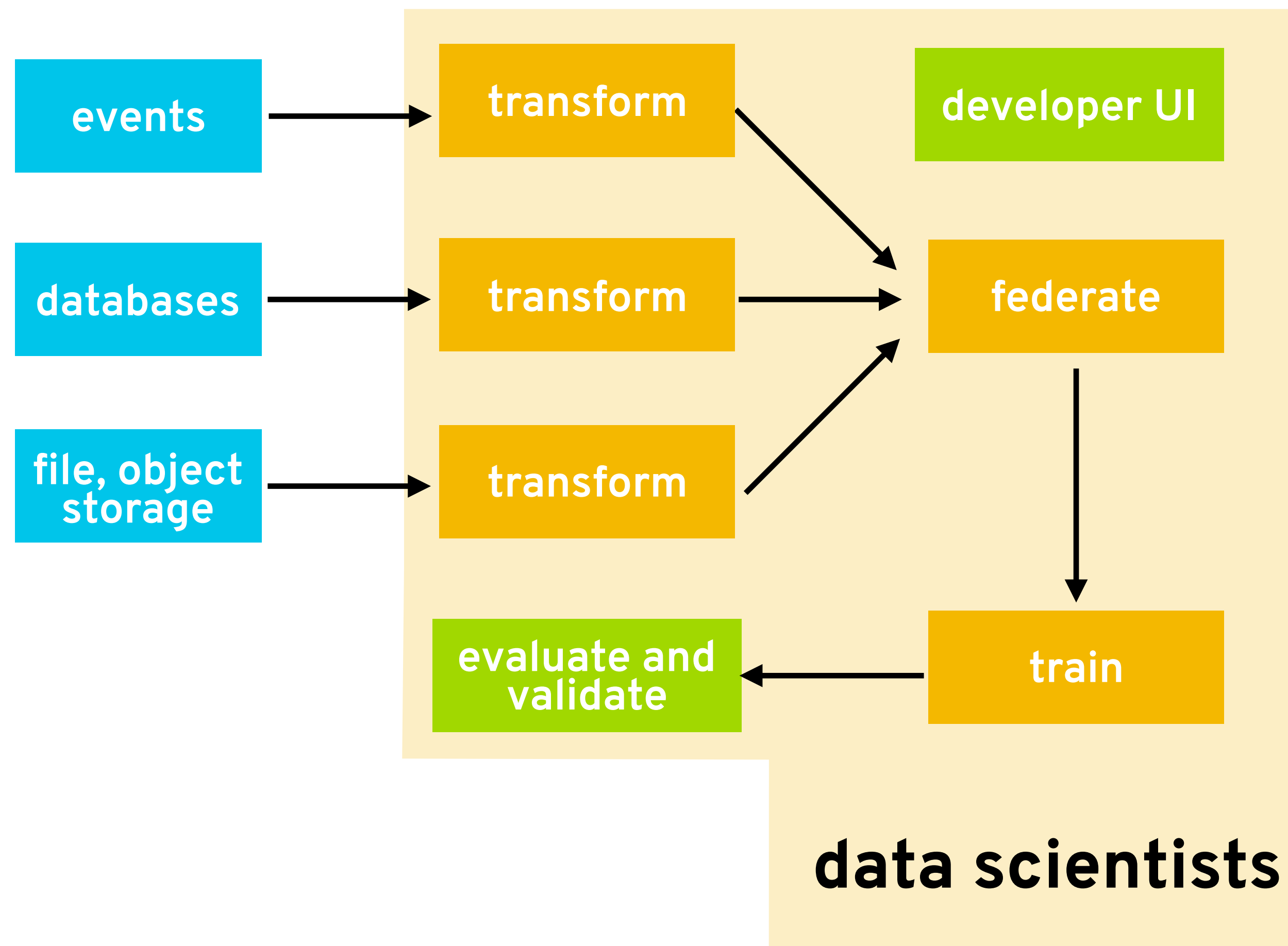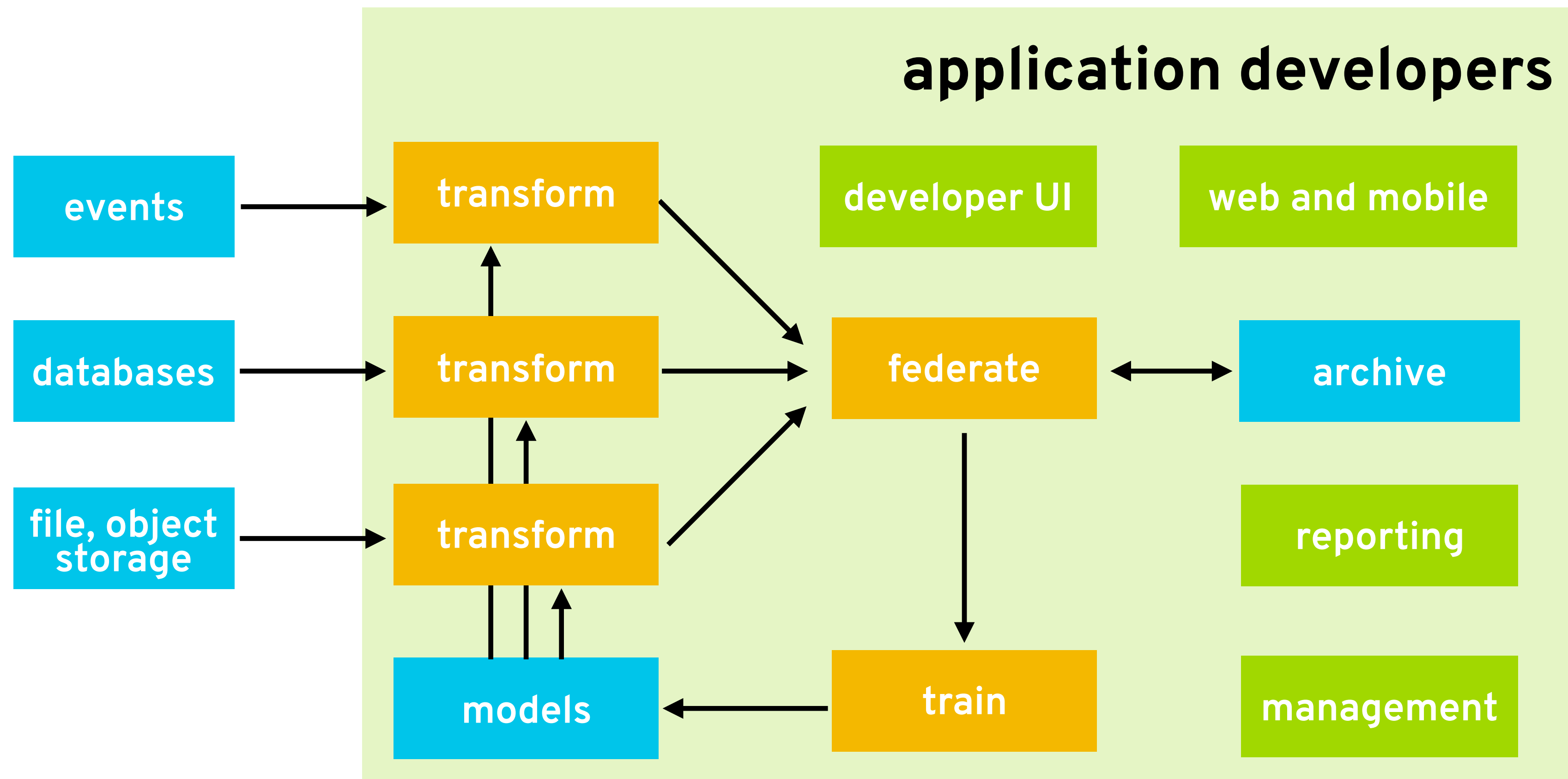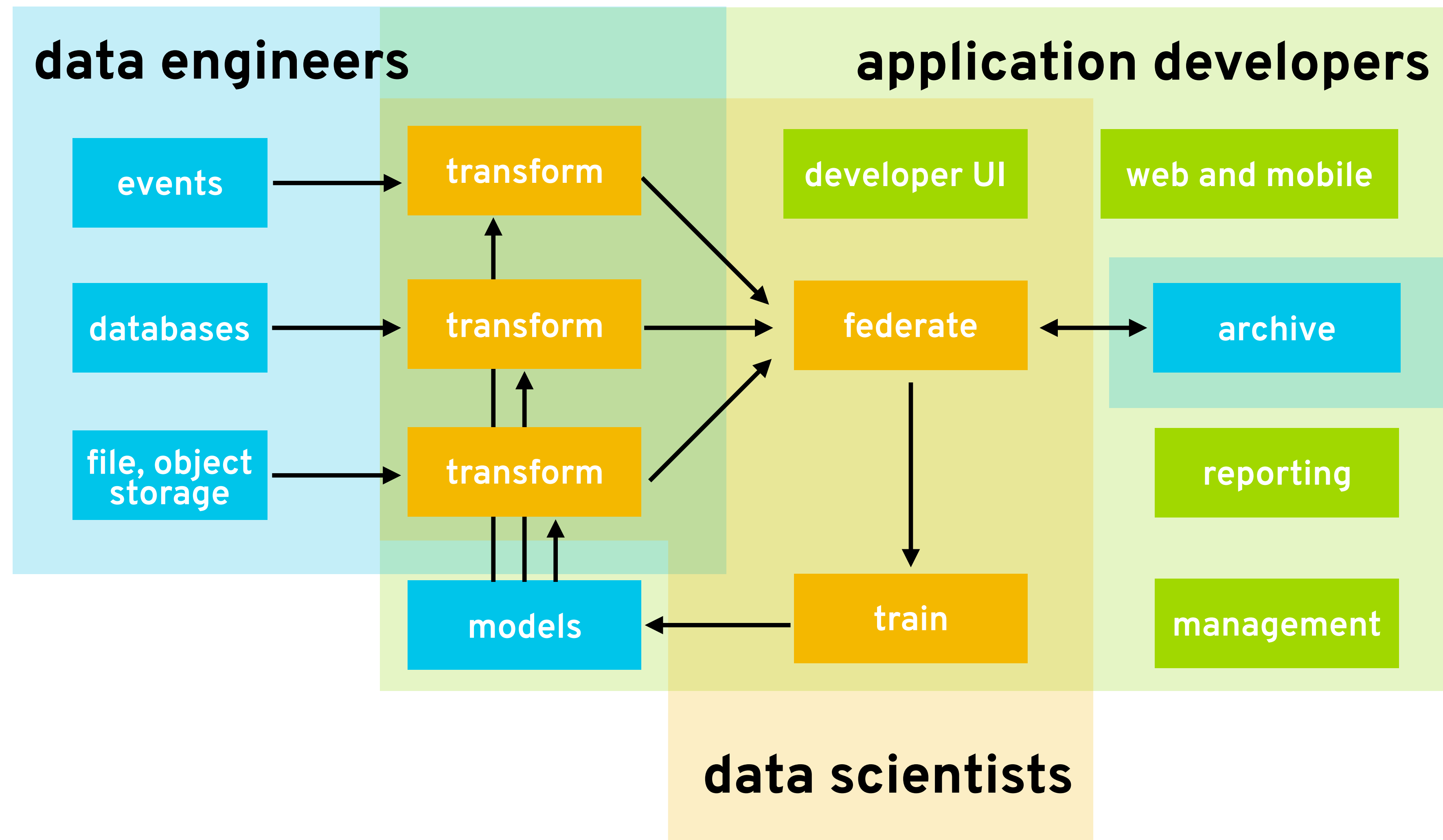
# Containers: great for application developers and data scientists

# What is a container?

%

```
% pip install numpy
```

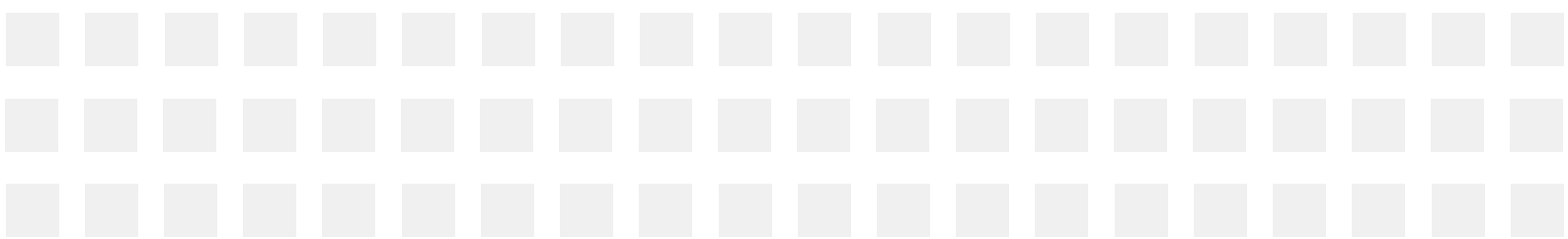| | |
|---|---|
| **executable** | `/usr/bin/pip` |
| **arguments** | `pip install numpy` |
| **environment** | `LANG=en_US USER=willb ...` |
| virtual memory | |
| file handles | |
| root filesystem | / |
| process table | |
| network routes | |

executable    /usr/bin/pip
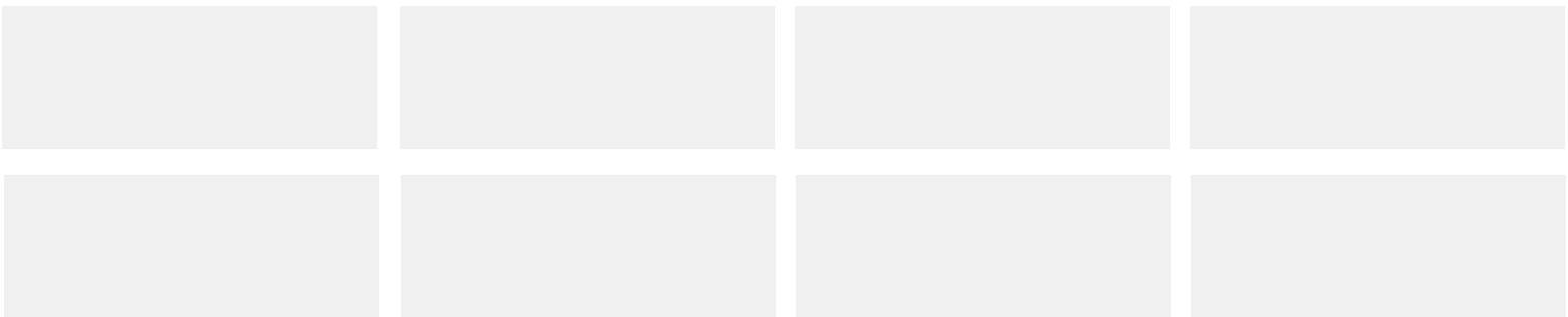
arguments     pip install numpy

environment   LANG=en_US USER=willb ...
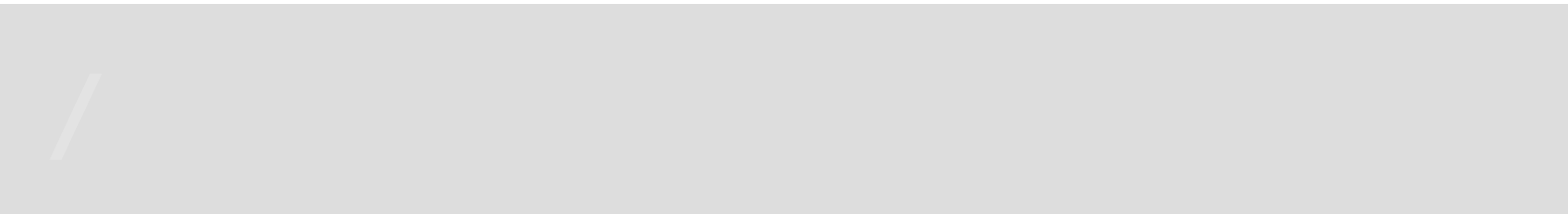
**virtual memory**

**file handles**

root filesystem    /

process table

network routes

executable    /usr/bin/pip

arguments    pip install numpy

Software Failure.   Press left mouse button to continue.
Guru Meditation #00000004.0000AAC0

root filesystem    /

process table

network routes

executable      /usr/bin/pip

arguments       pip install numpy

```
Software Failure.   Press left mouse button to continue.
         Guru Meditation #00000004.0000AAC0
```
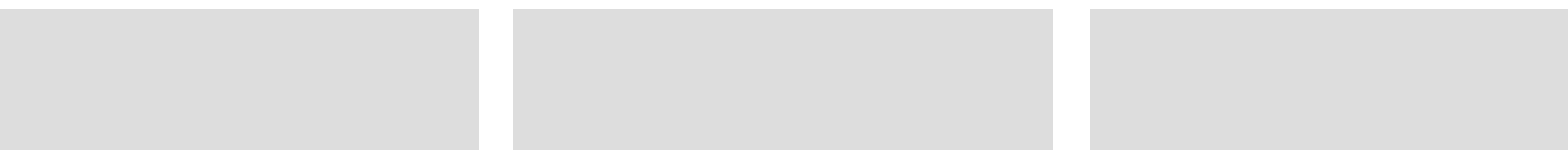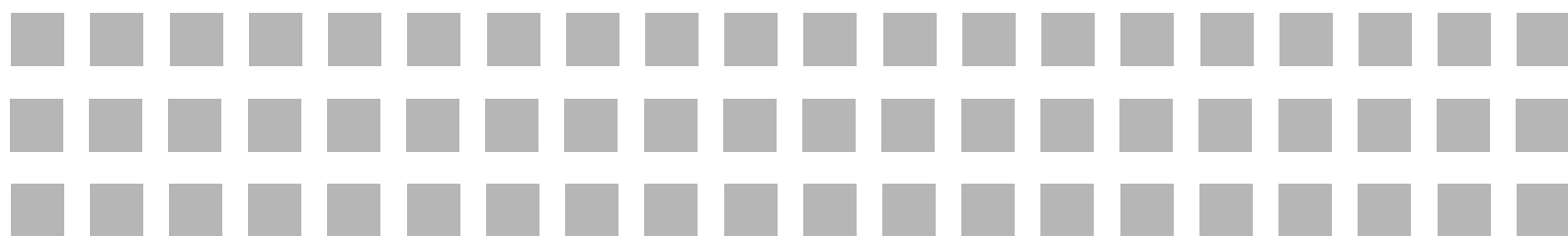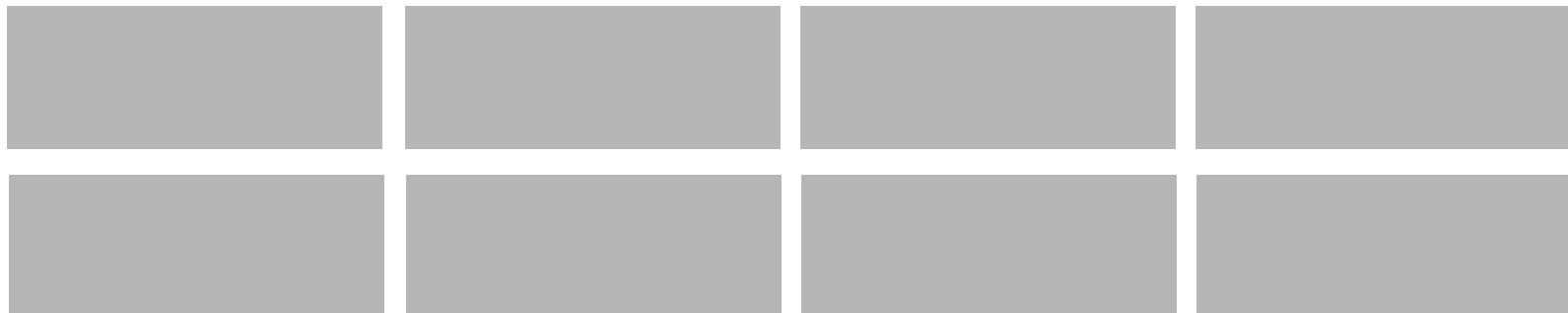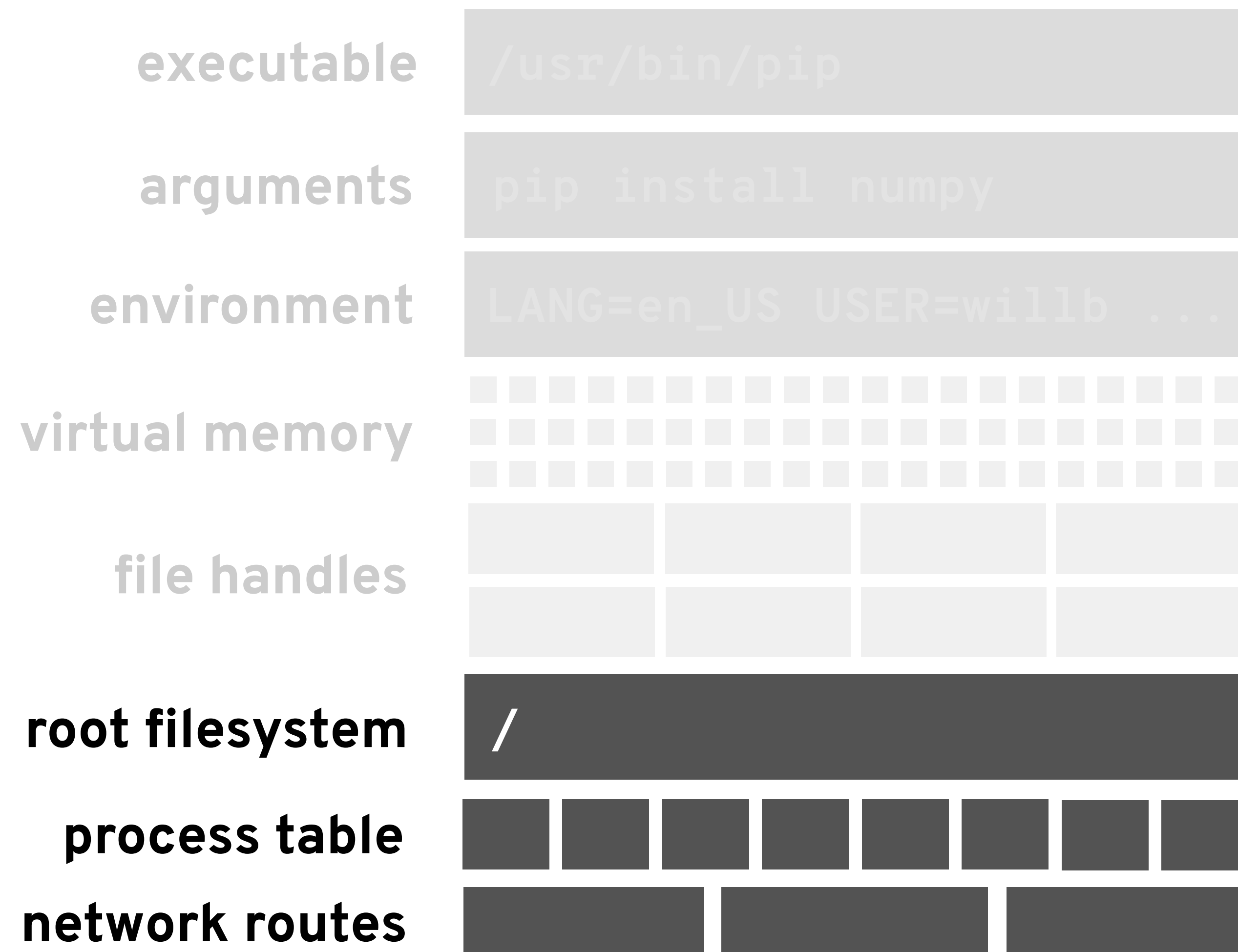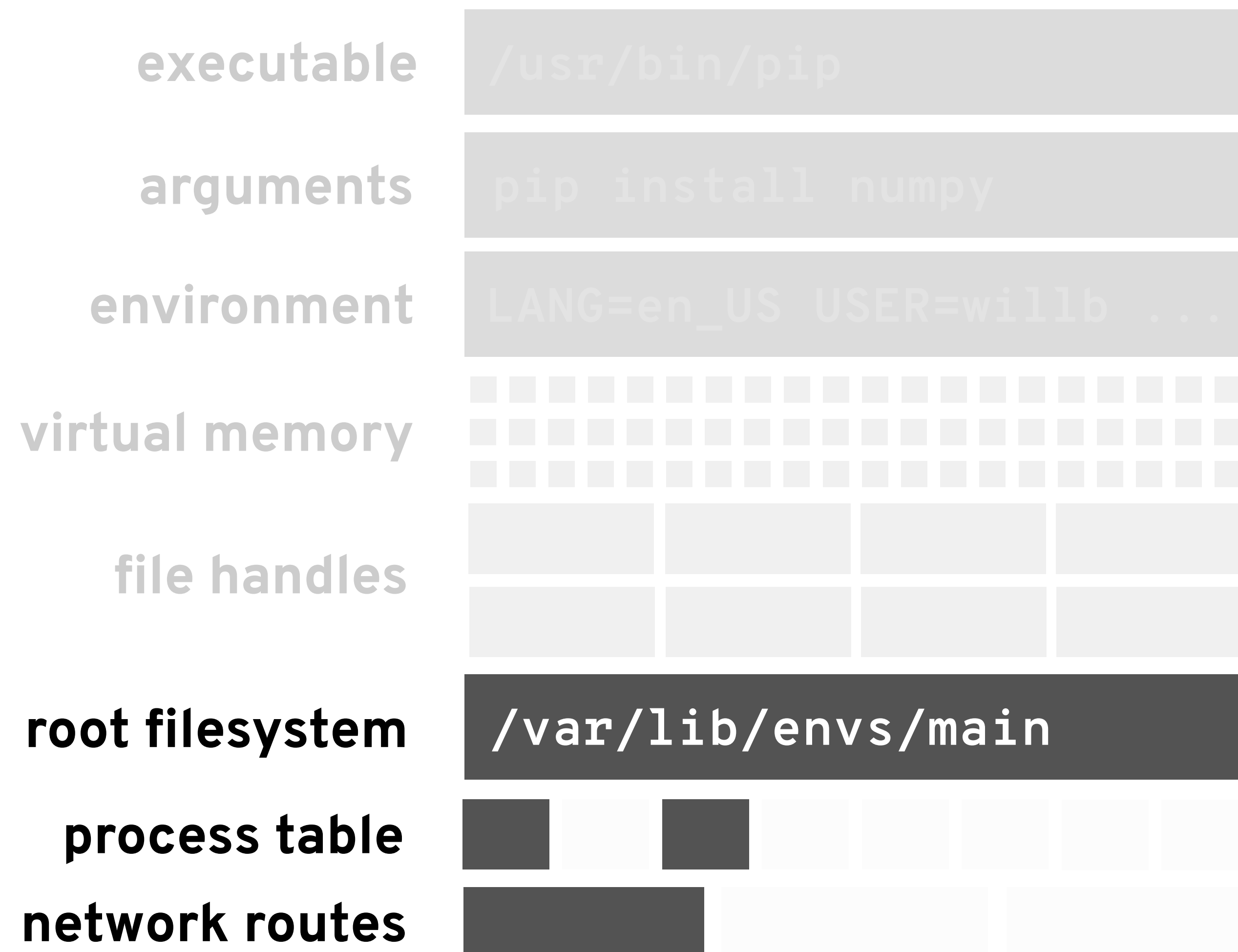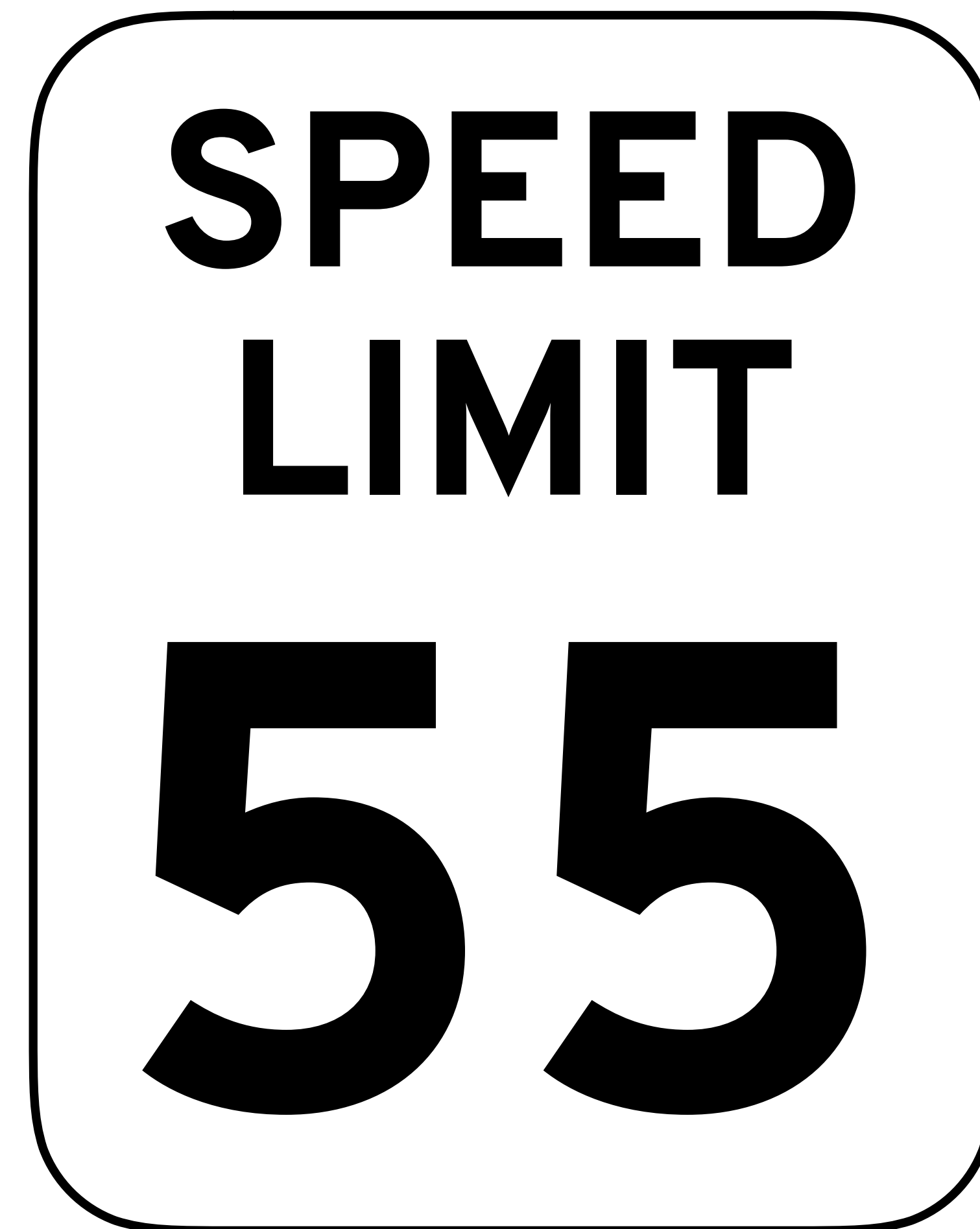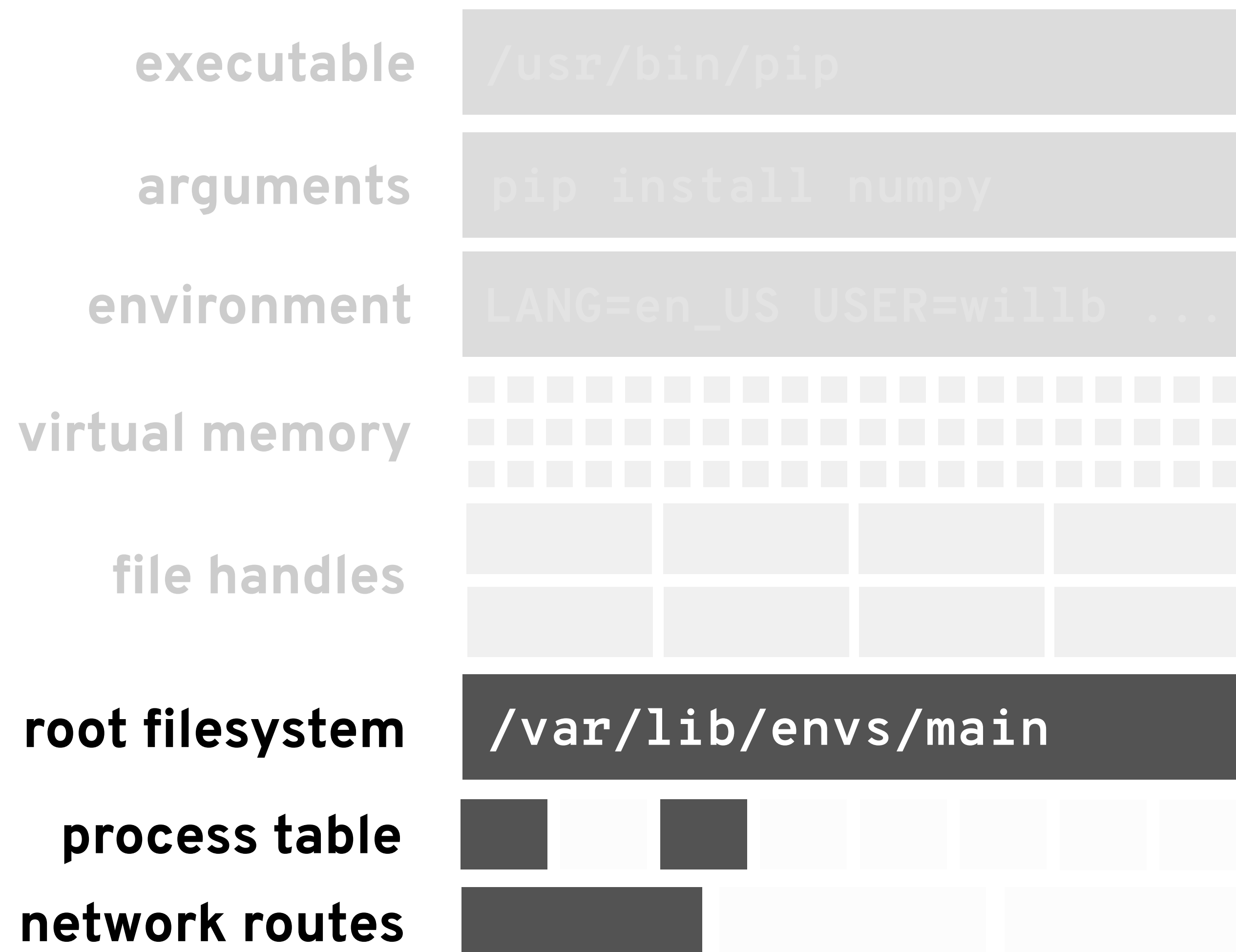
root filesystem    /

process table

network routes

executable    /usr/bin/pip

arguments    pip install numpy

environment    LANG=en_US USER=willb ...

virtual memory

file handles

**root filesystem**    /

**process table**

**network routes**

executable    /usr/bin/pip

arguments    pip install numpy

environment    LANG=en_US USER=willb ...

virtual memory

file handles

**root filesystem**    `/var/lib/envs/main`

**process table**

**network routes**

| | |
|---|---|
| executable | /usr/bin/pip |
| arguments | pip install numpy |
| environment | LANG=en_US USER=willb ... |
| virtual memory | |
| file handles | |
| **root filesystem** | /var/lib/envs/main |
| **process table** | |
| **network routes** | |

**SPEED LIMIT**

**55**

| | | |
|---|---|---|
| executable | /usr/bin/pip | /usr/bin/pip |
| **arguments** | `pip install numpy` | `pip install riskylib` |
| environment | LANG=en_US USER=willb ... | LANG=en_US USER=willb ... |
| virtual memory | | |
| file handles | | |
| **root filesystem** | `/var/lib/envs/main` | `/var/lib/envs/risky` |
| **process table** | | |
| **network routes** | | |

# Immutable images

base image

# Immutable images

configuration and installation recipes

base image

# Immutable images

| | |
|---|---|
| **user application code** | `a6afd91e`<br>`6b8cad3e` |
| **configuration and installation recipes** | `33721112`<br>`e8cae4f6`<br>`2bb6ab16`<br>`a8296f7e` |
| **base image** | `979229b9` |

# Stateless microservices

# Stateless microservices

# Stateless microservices
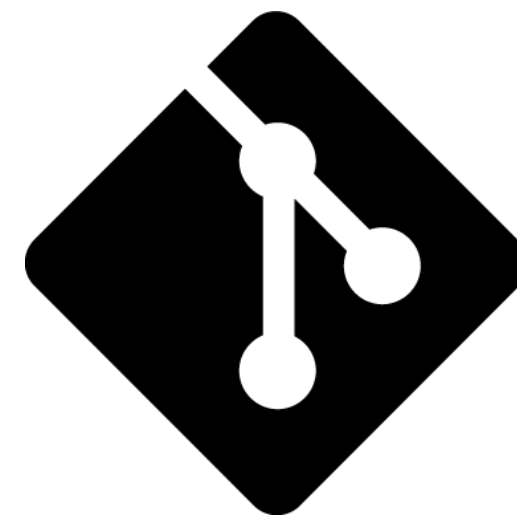
# Stateless microservices
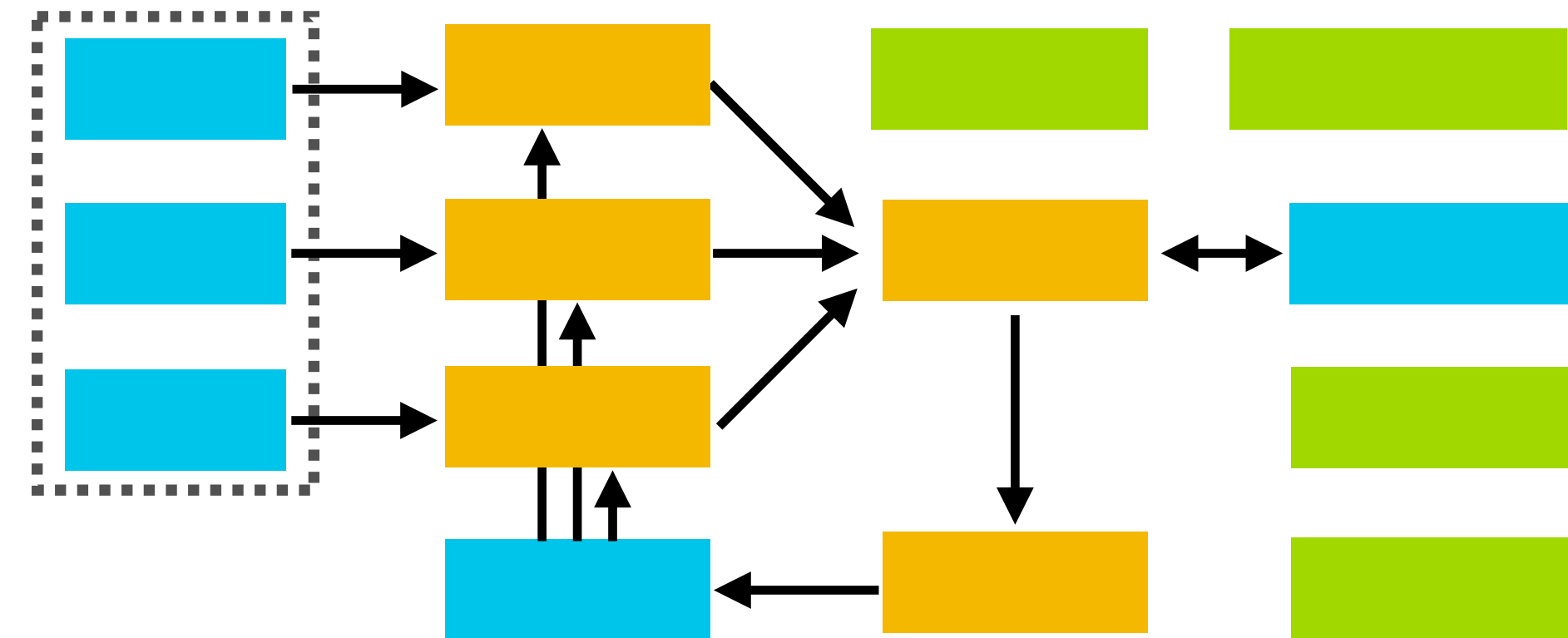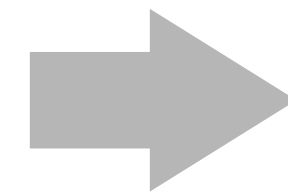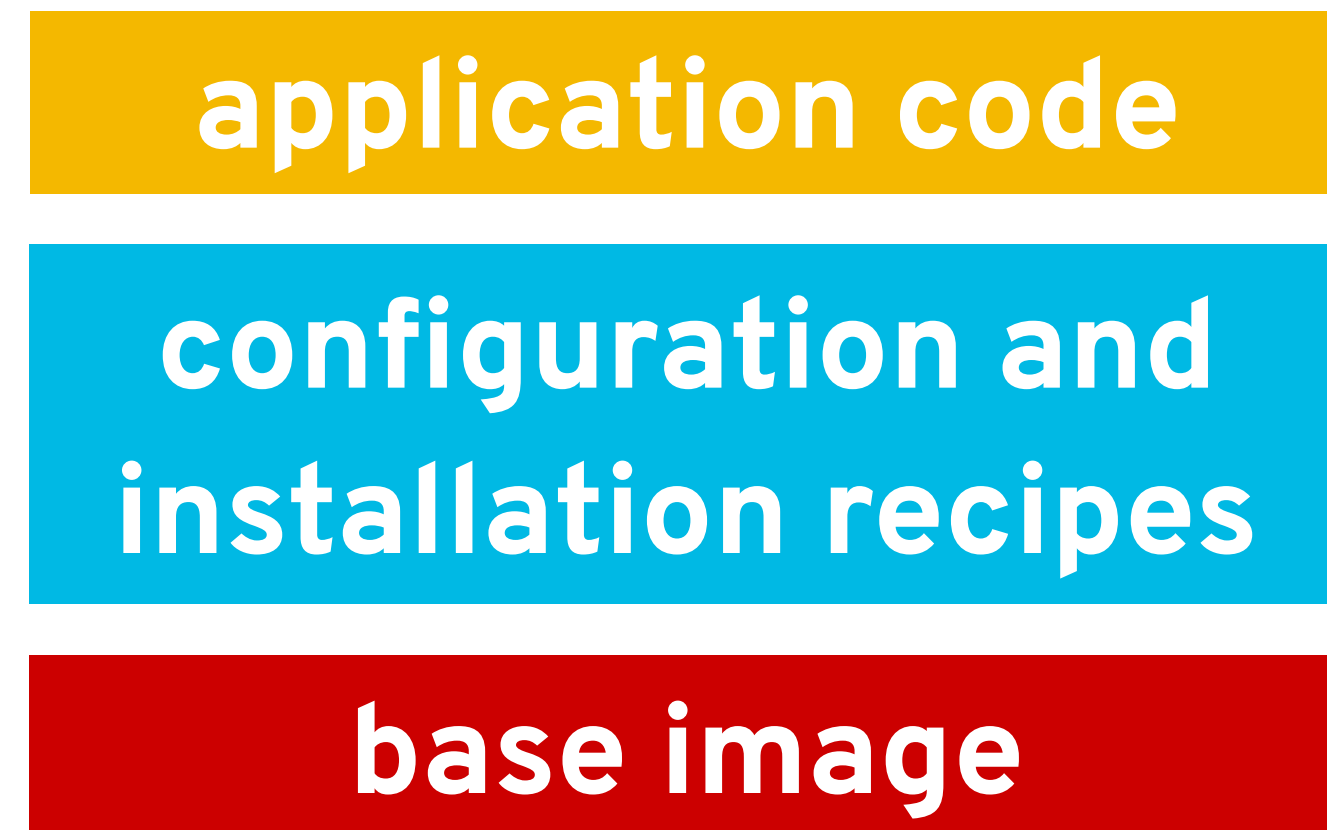
# Declarative app configuration

# Integration and deployment

# Integration and deployment

# Integration and deployment

# Integration and deployment



OK!

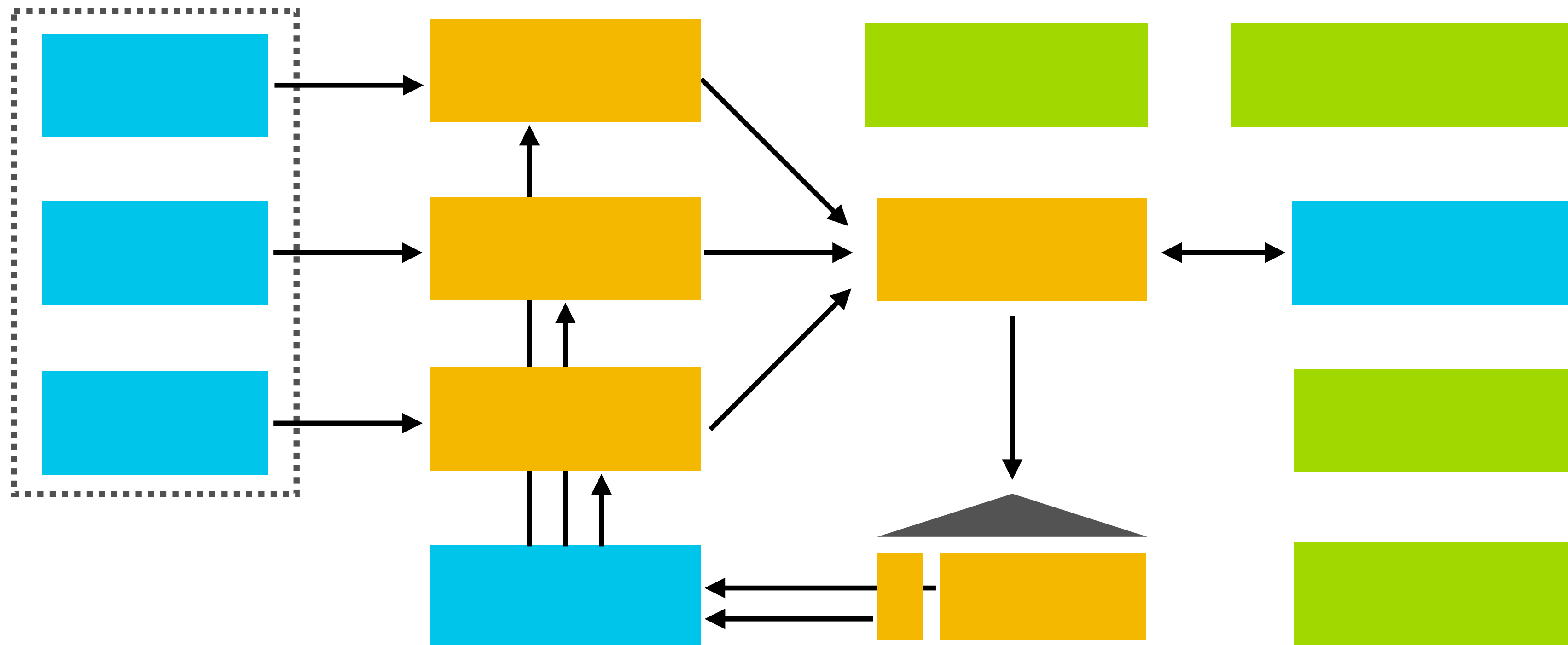application code

configuration and installation recipes

base image

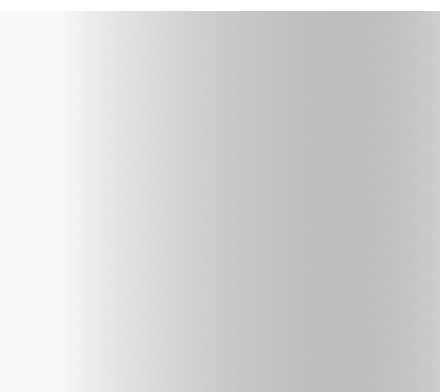# Integration and deployment

# Flexible service routing

# Flexible service routing

# How can we build *intelligent applications* in containers?
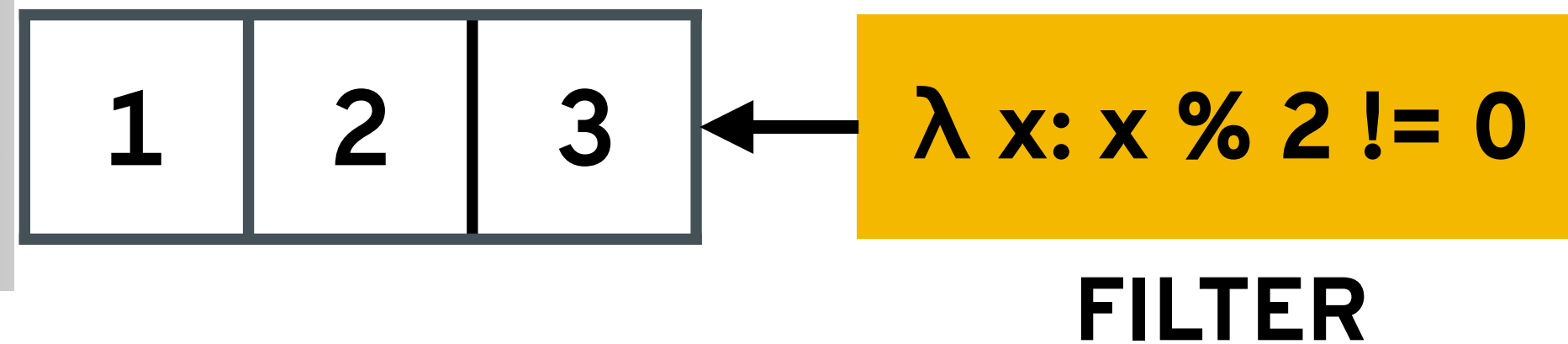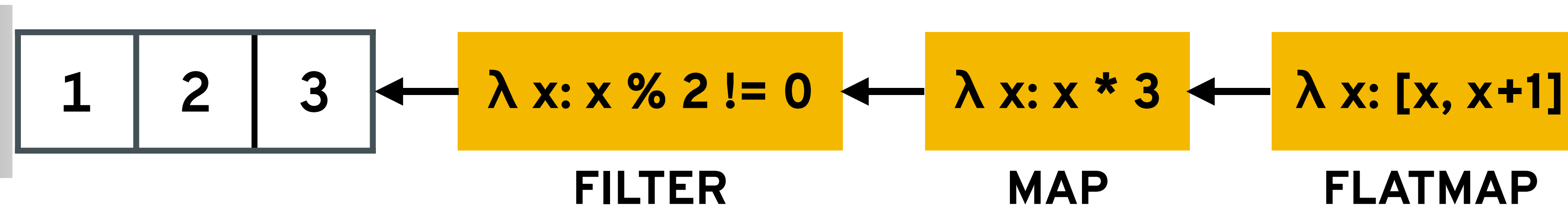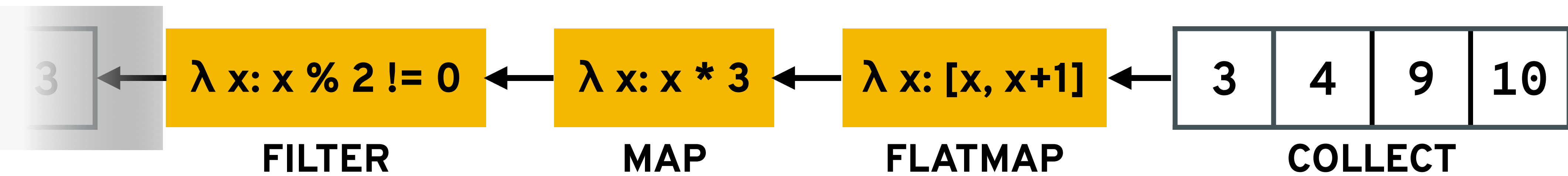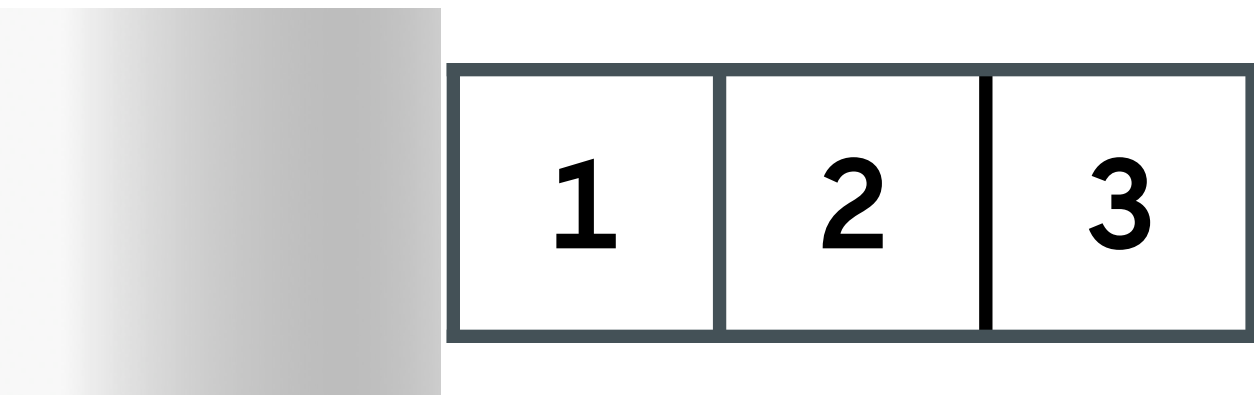
# Apache Spark

# Apache Spark

| 1 | 2 | 3 |
|---|---|---|

# Apache Spark

# Apache Spark

| 1 | 2 | 3 |
|---|---|---|

← **λ x: x % 2 != 0**

**FILTER**

← **λ x: x * 3**

**MAP**

# Apache Spark

| 1 | 2 | 3 | ← | λ x: x % 2 != 0 | ← | λ x: x * 3 | ← | λ x: [x, x+1] |

**FILTER**  **MAP**  **FLATMAP**

| 1 | 2 | 3 |

1 | 2 | 3 ← λ x: x % 2 != 0

**FILTER**

3 ← **λ x: x % 2 != 0**

**FILTER**

**3**

**λ x: x % 2 != 0** ← **λ x: x * 3** ← **λ x: [x, x+1]** ← **CACHE** ← | 3 | 4 | 9 | 10 |

**FILTER**      **MAP**      **FLATMAP**      **SAVE AS TEXT FILE**

# Spark core

| 1 | 2 | 3 | | 4 | 5 | 6 | | 7 | 8 | 9 | | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|----|----|----|

| executor | executor | executor | executor | master |

| 2 | 4 | 6 | | 8 | 10 | 12 | | 14 | 16 | 18 | | 20 | 22 | 24 | | λ x: x * 2 |

| executor | executor | executor | executor | master |

| λ x: x * 2 | λ x: x * 2 | λ x: x * 2 | λ x: x * 2 |

| 2 | 4 | 6 |
|---|---|---|

| 8 | 10 | 12 |
|---|---|---|

executor

executor

λ x: x * 2

λ x: x * 2

| 20 | 22 | 24 |
|----|----|----|

λ x: x * 2

executor

master

λ x: x * 2

| 2 | 4 | 6 | | 8 | 10 | 12 | | 14 | 16 | 18 | | 20 | 22 | 24 | | λ x: x * 2 |

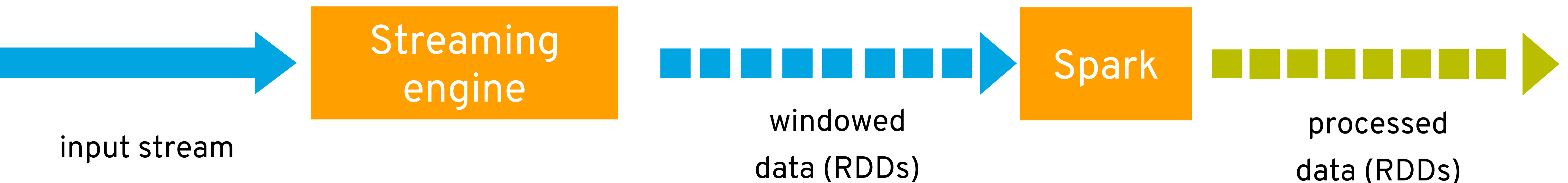| executor | executor | executor | executor | master |

| λ x: x * 2 | λ x: x * 2 | λ x: x * 2 | λ x: x * 2 |

# Streaming data

Goal: use the **same abstraction** for batch and "streaming" (micro-batch) data by **dividing a stream** into **many small RDDs.**

input stream

# Streaming data

Goal: use the **same abstraction** for batch and "streaming" (micro-batch) data by **dividing a stream** into **many small RDDs.**

Streaming engine

input stream

# Streaming data

Goal: use the **same abstraction** for batch and "streaming" (micro-batch) data by **dividing a stream** into **many small RDDs.**

input stream → Streaming engine → windowed data (RDDs)

# Streaming data

Goal: use the **same abstraction** for batch and "streaming" (micro-batch) data by **dividing a stream** into **many small RDDs.**

input stream → | Streaming engine | ┈┈┈> windowed data (RDDs) | Spark |

# Streaming data

Goal: use the **same abstraction** for batch and "streaming" (micro-batch) data by **dividing a stream** into **many small RDDs.**

Streaming engine

input stream

windowed data (RDDs)

Spark

processed data (RDDs)

# Structured queries

The capacity to run arbitrary code in RDDs is powerful but comes with **an important tradeoff**: Spark can't rearrange RDD programs to improve their performance.

Writing Spark programs with a **query DSL** allows Spark to generate **optimized execution plans**.

# Structured query in Spark

**SQL interface** (unchecked syntax or semantics)
```
SELECT word, COUNT(*) FROM words GROUP BY word
```

**Data frame interface** (semantics checked at run-time)
```
words.groupBy('word').count()
```

**Dataset interface** (mostly checked at compile-time)

# Query planning

```
SELECT * FROM A, B WHERE
  A.ID = B.ID AND
  uncommon(A.X) AND
  extremelyRare(B.Y)
```

# A naïve plan

# A naïve plan

JOIN

# A naïve plan

JOIN

# A naïve plan

# A naïve plan



**FILTER**

# A naïve plan

**FILTER**

# An optimized plan

# An optimized plan

**FILTER**          **FILTER**

# An optimized plan



**FILTER**        **FILTER**

# An optimized plan



JOIN

# An optimized plan

Structured streaming combines **stream processing** with **query planning** for **high-performance analytics on events**!

# The Kappa architecture

events

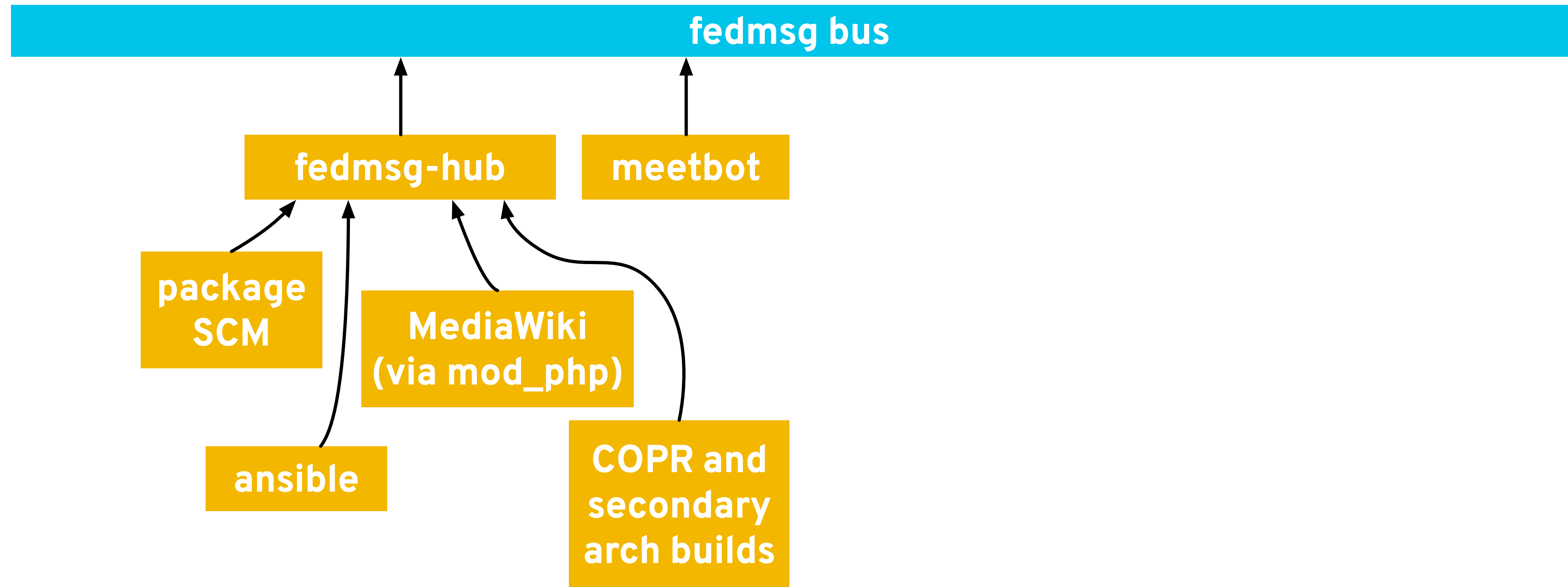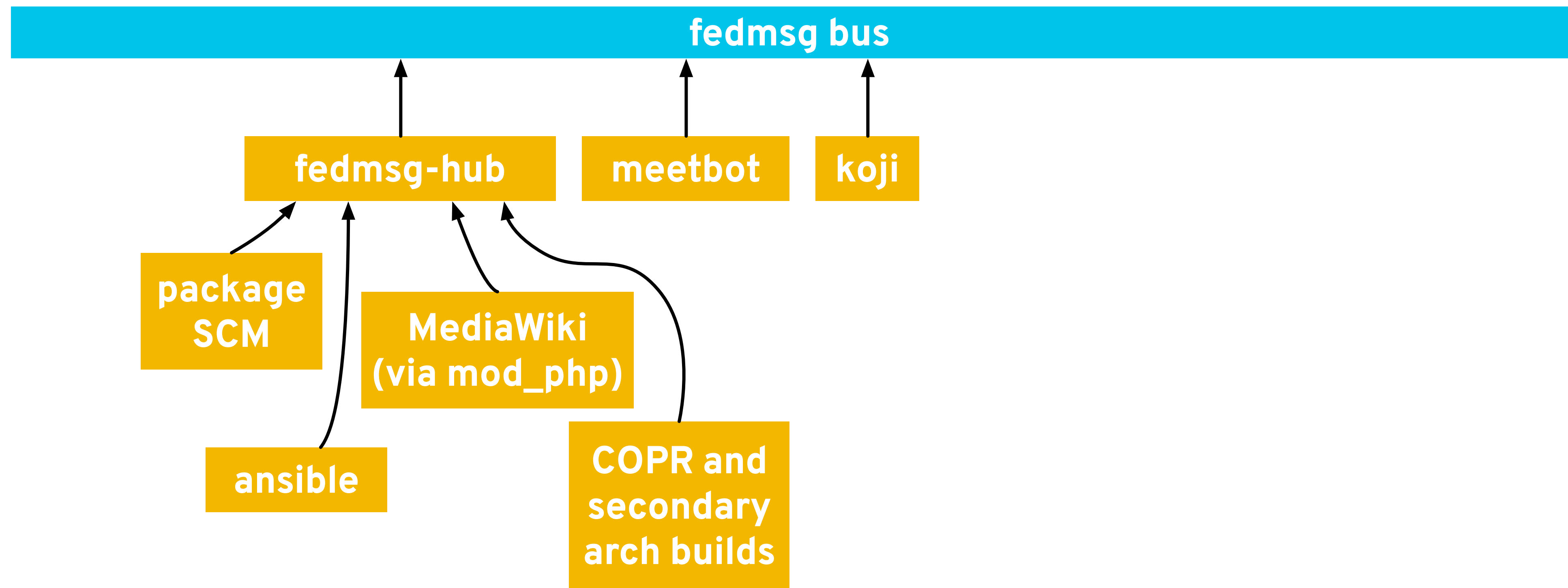# The Kappa architecture

# The Kappa architecture

events → queue for "raw data" topic

# The Kappa architecture

events → queue for "raw data" topic

transform

# The Kappa architecture
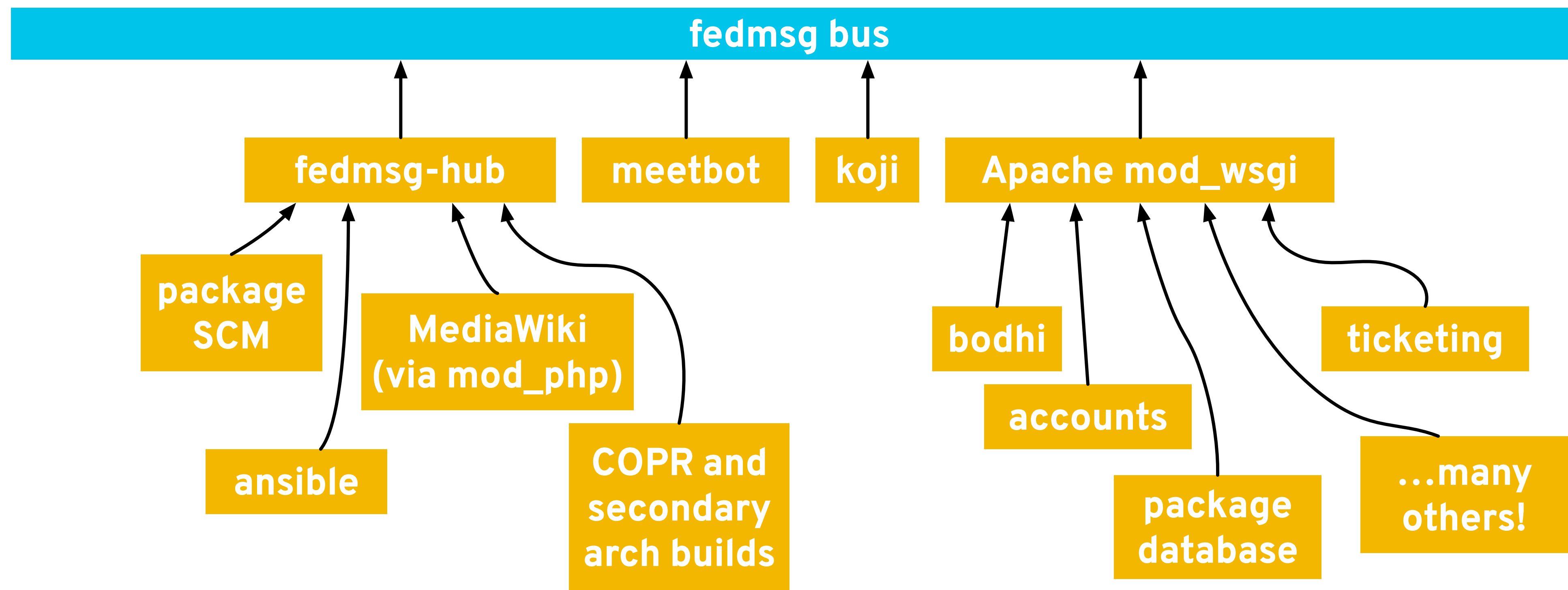
# The Kappa architecture

# The Kappa architecture

# The Kappa architecture

# The Kappa architecture

# fedmsg bus

# Apache Kafka

# Apache Kafka

# Apache Kafka

| 1 | 3 | 7 | 11 |
|---|---|---|----|

| 2 | 4 | 5 | 10 |
|---|---|---|----|

| 6 | 8 | 9 | 12 |
|---|---|---|----|

# Multitenant compute clusters

**Resource manager**

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Spark executor

Shared FS / object store

Databases

# Multitenant compute clusters



**Resource manager**

Spark executor
Spark executor
Spark executor
Spark executor
Spark executor
Spark executor

Shared FS / object store

Databases

# One cluster per application

# One cluster per application

# radanalytics.io

An **open-source community** enabling **intelligent applications** on **OpenShift**

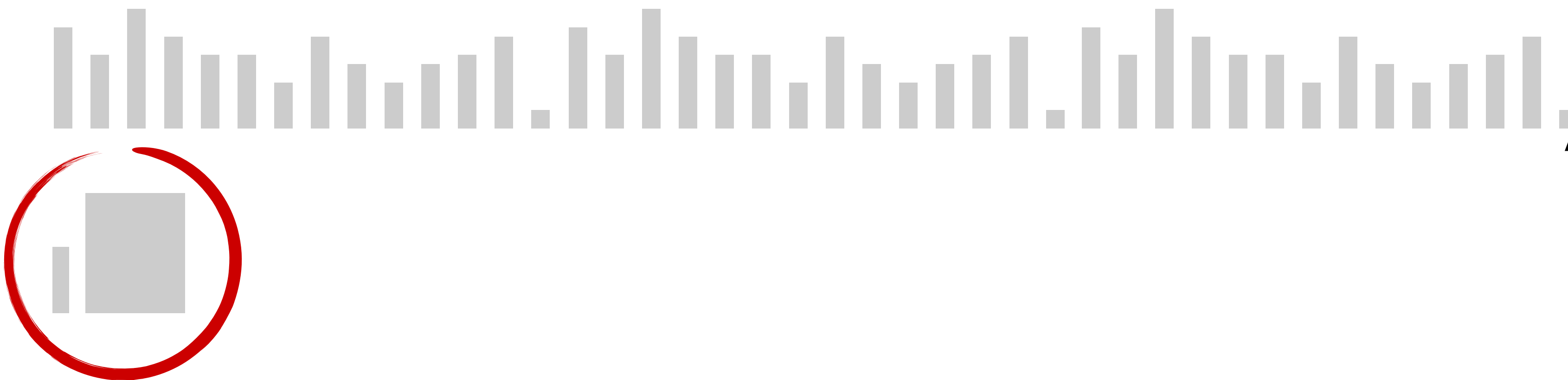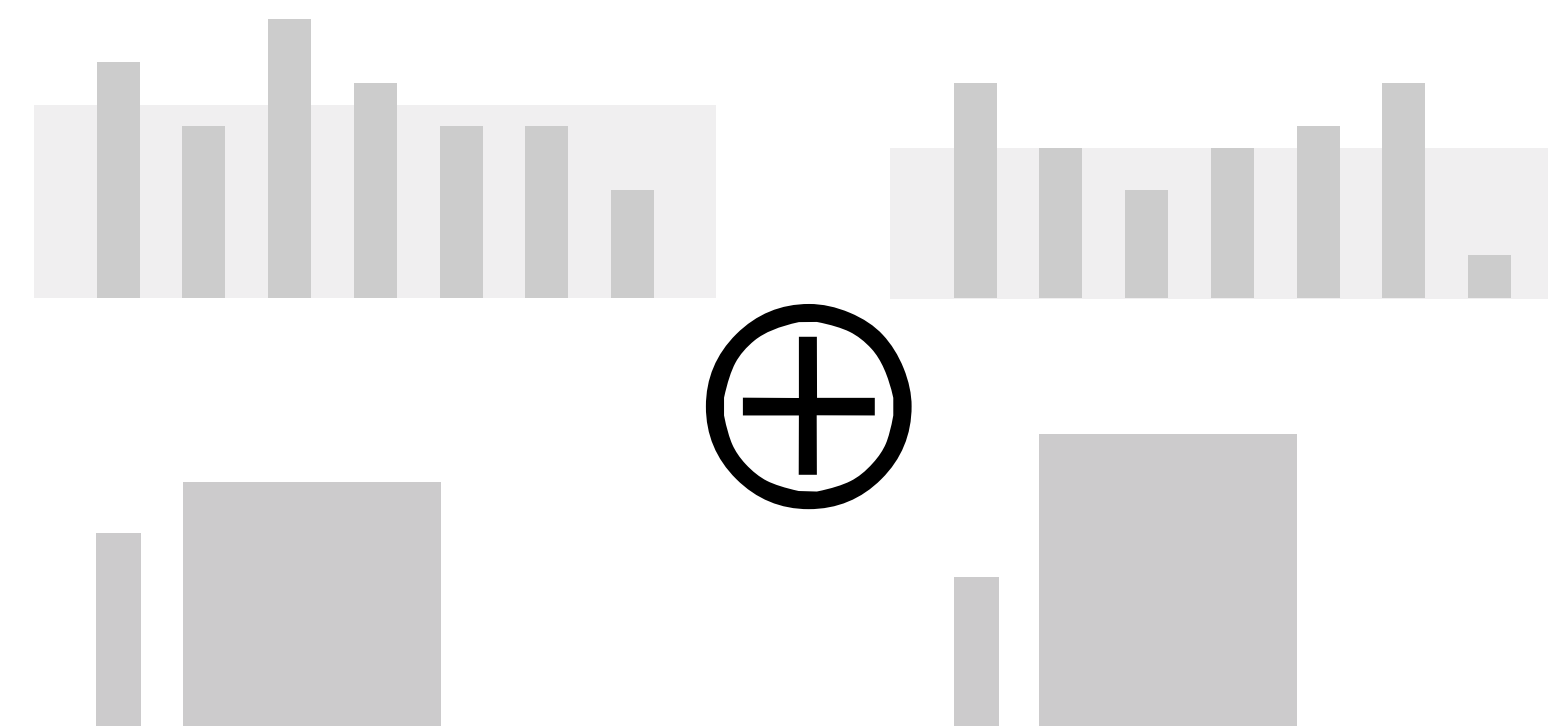Tooling to manage **Apache Spark**, **Jupyter notebooks**, and **TensorFlow training** and **model serving**

Interactivity time!

# http://bit.ly/streaming-bds18
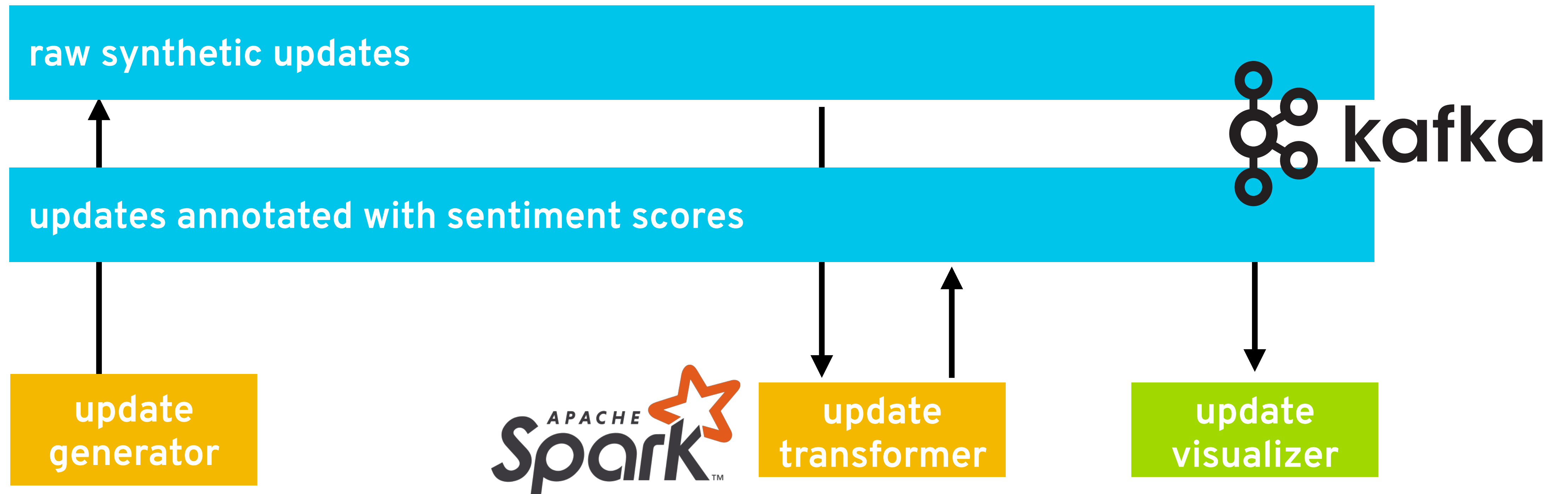
# Searching for solutions

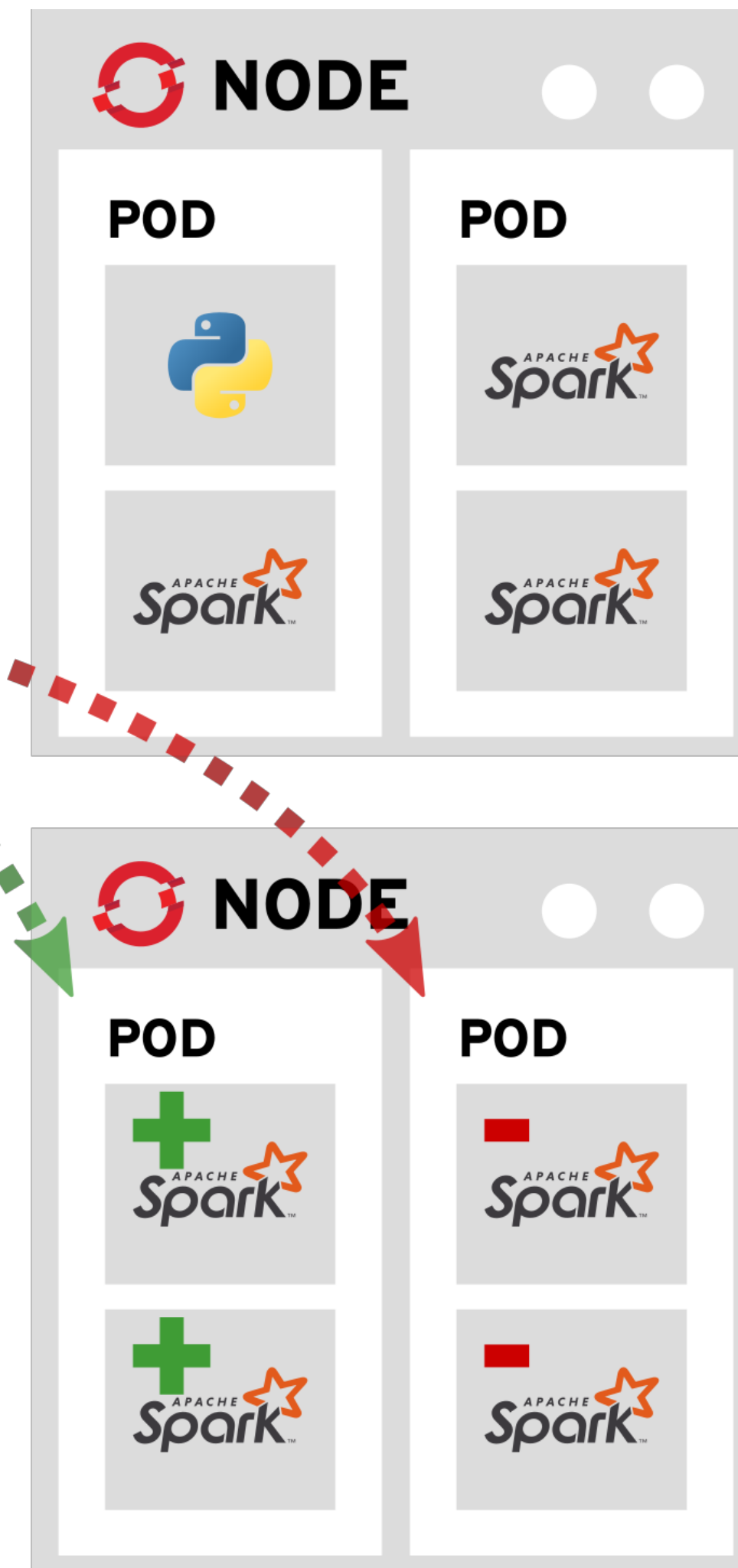# Consider possible interactions

# Plan for common patterns

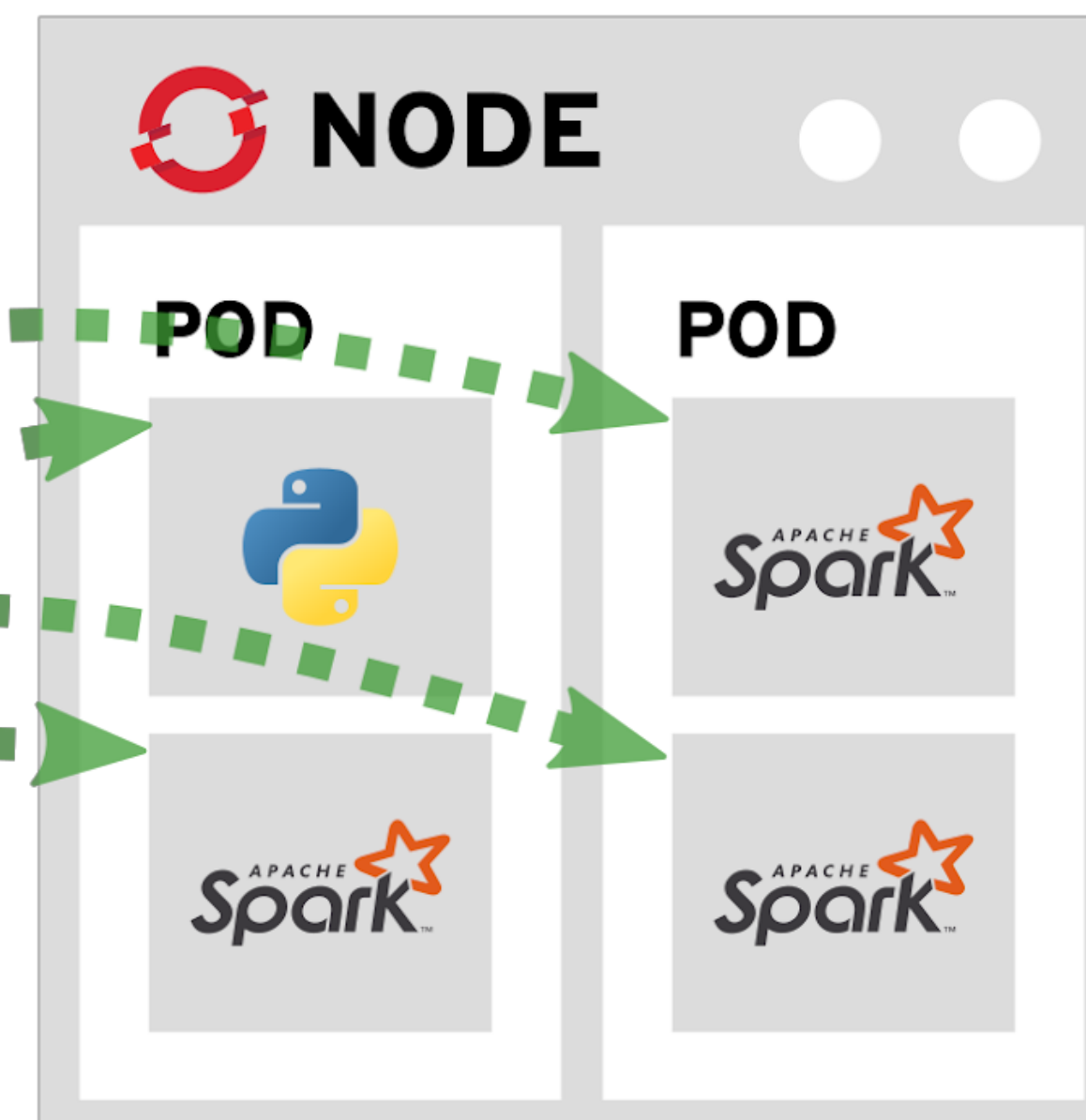Ingest > Process > Publish

oshinko webui

**oshinko source-to-image**