

CC5212-1

PROCESAMIENTO MASIVO DE DATOS

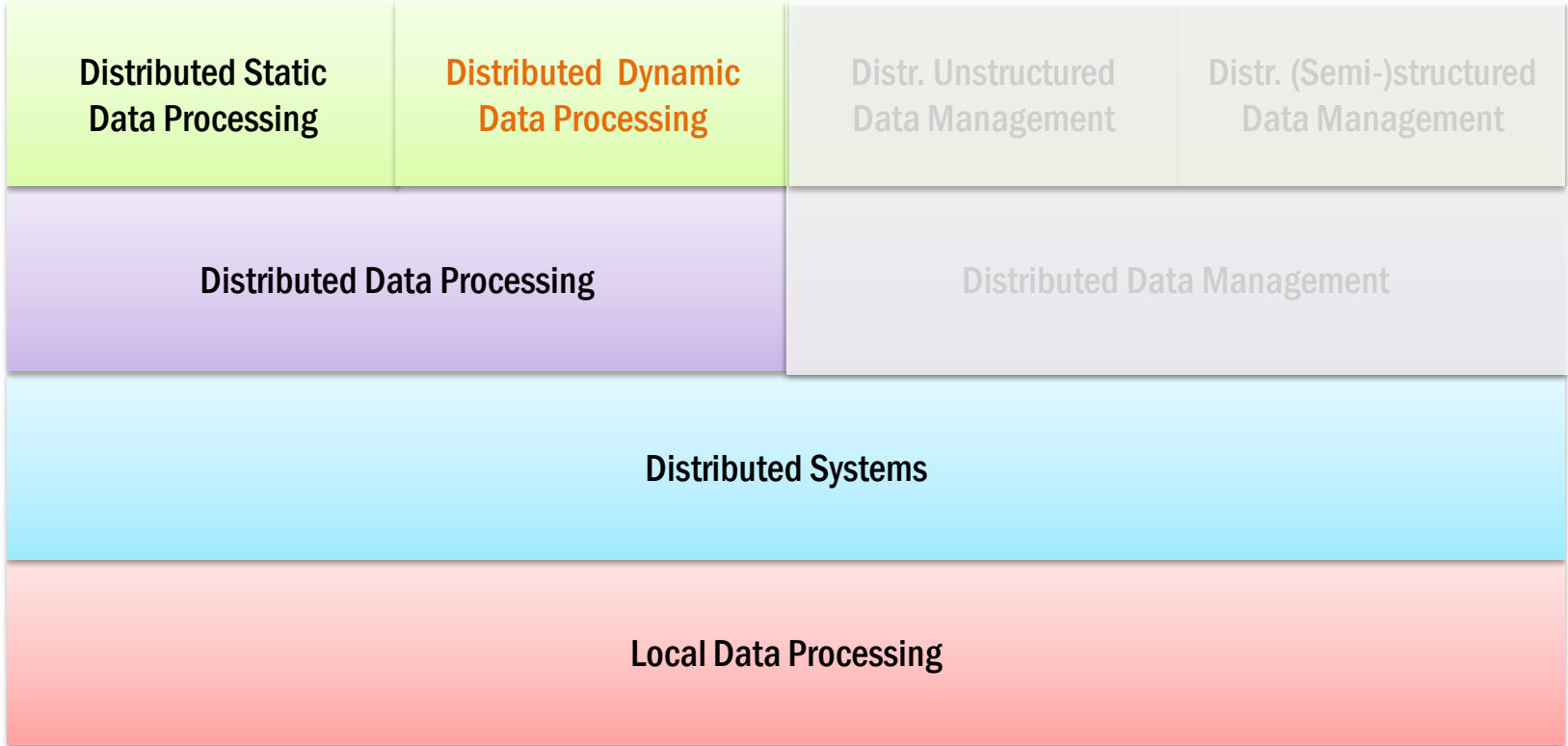
OTOÑO 2023

Lecture 6

Streaming: Kafka

Aidan Hogan

aidhog@gmail.com

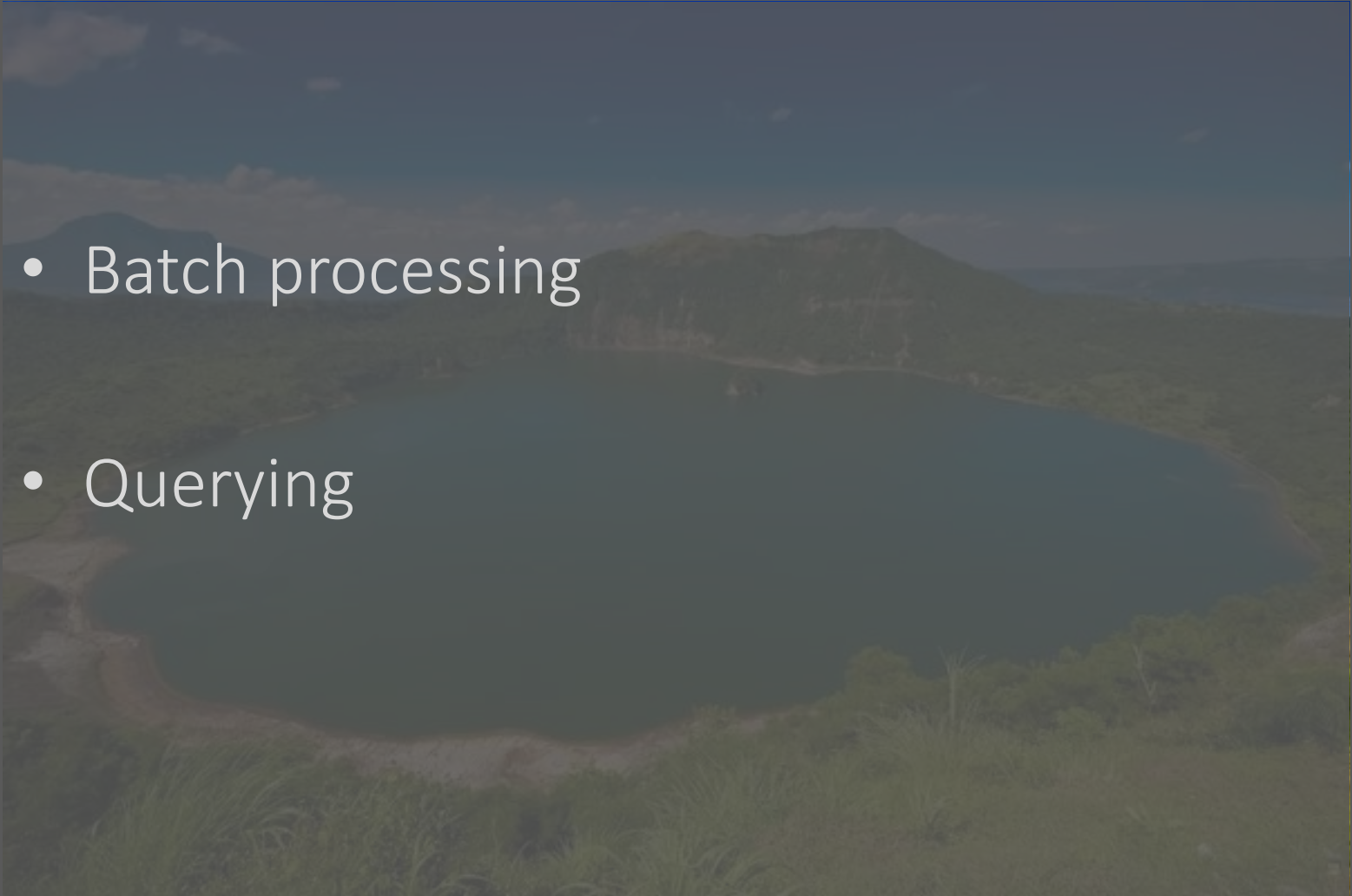


Files



Files

- Batch processing
- Querying

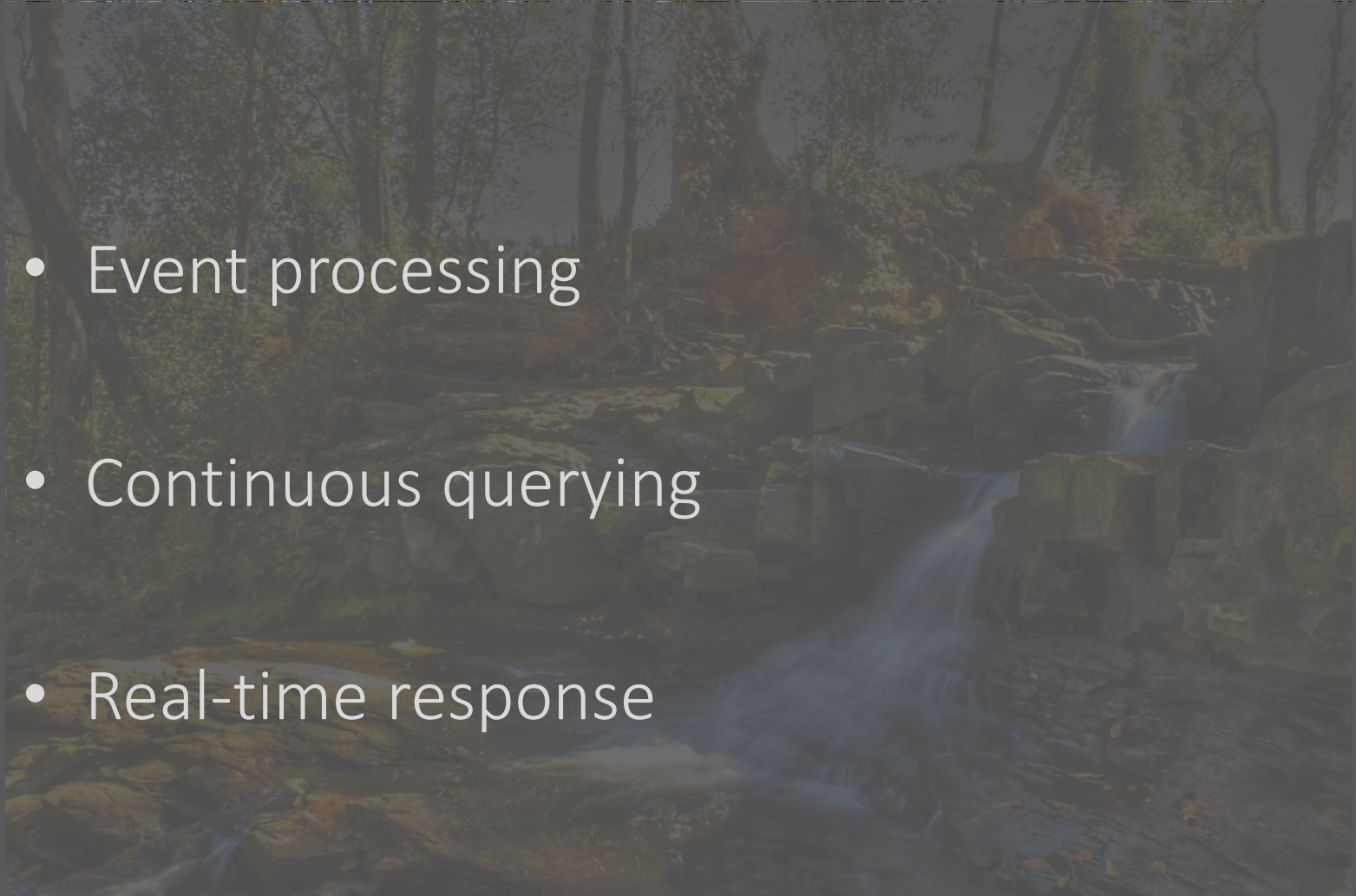


Streams



Streams

- Event processing
- Continuous querying
- Real-time response



Applications: Social Media Analytics



Applications: Social Media Analytics

- Event processing
 - Kitten video goes viral
 - Burst of tweets about earthquakes
- Continuous querying
 - Track sentiment for company's products
 - Monitor popular users tweeting about me
- Real-time response
 - Put Emergency Services on alert
 - Schedule Quality Control (QC) review

Applications: Log Monitoring

```
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Cache-Control' => 'no-cache' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Connection' => 'keep-alive' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Pragma' => 'no-cache' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Accept' => 'application/json, text/javascript, */*; q=0.01' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Accept-Encoding' => 'gzip, deflate' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Accept-Language' => 'en-gb,en;q=0.5' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Host' => 'vx-garcia.wcn.co.uk' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Referer' => 'http://vx-garcia.wcn.co.uk/vx/lang-en-GB/config-jail/channel-1
d27dad84/wid-4/ats/recruiter/profile/edit' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'User-Agent' => 'Mozilla/5.0 (X11; Linux x86_64; rv:17.0) Gecko/17.0 Firefox
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Content-Length' => '134' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Content-Type' => 'application/x-www-form-urlencoded; charset=UTF-8' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'Cookie' => 'utma=1.893613000.1353586388.1356003545.1356012074.105; utmz=1355835790.3; utmc=164338489.1362821095.1354014424.1354641679.1355835790.3; utmcsr=google|utmccn=(organ
|utmccn=(direct)|utmcmd=(none); utma=164338489.1362821095.1354014424.1354641679.1355835790.3; utmz=164338489.1355835790.3.3.utmcsr=google|utmccn=(organ
provided); wcn_ats_session=000097a096228f7b2bdafa54194513879815e69871128c680440ee76ee108d9d39d6c44ddd261dd4ffa; utmc=1; su_user=0; utmb=1.89.9.135601659
0ec1126768805e55e2137f801ea5d9ad42d9f35dc1f31f70a568b822e61ace6f080f6' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 158] 'X-Requested-With' => 'XMLHttpRequest' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 4] [WCN.Role.PathMunge 173] In parse_url() - user requested access to 'vx-garcia.wcn.co.uk', 'lang-en-GB/config-jail
-1/xf-cbf4d27dad84/wid-4/ats/recruiter/profile/map_update/1'
[Thu Dec 20 15:41:01 2012] [notice] Apache/2.2.16 (Debian) mod_perl/2.0.4 Perl/v5.10.1 configured -- resuming normal operations
[INFO] [47152f7f] [2012/12/20 15:41:01 48] [WCN.Role.PathMunge 747] c->set_system: 51 (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.DBIC 388] Going to set system database to system '51', jail '1'
[INFO] [47152f7f] [2012/12/20 15:41:01 36] [WCN.Role.PathMunge 718] Setting brand to be '2'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PathMunge 791] Set current language to 'en-GB'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 2] [WCN.Role.PathMunge 319] Cookie for 'recruiter' => '97a096228f7b2bdafa54194513879815e6987112
[DEBUG] [47152f7f] [2012/12/20 15:41:01 6] [WCN.Role.PerformanceLogger 175] *** Request Params (4) ***
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 189] '_vxXSRF_Token' => '4ddc5de6bed3d3760b59c430cbeae48a6f0c8e95' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 189] 'code version' => '1355995679' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 189] 'submitted_via_ajax' => 'true' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 189] 'datafield_53274_1_1[]' => '1798' (1)
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 193] *** Uploads (0) ***
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 215] PERFORMANCE: prepare took 0.11599937057495 secs
[DEBUG] [47152f7f] [2012/12/20 15:41:01 1] [WCN.Role.Session 142] User's session id is '97a096228f7b2bdafa54194513879815e6987112', on server '0'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 17] [WCN.AccessControl 233] Going to _cache_role_profile_rules for recruiter '1', role profile '20'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 16] [WCN.Controller.Root 64] ** Enter root auto
[INFO] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Controller.Root 65] ** User requested access to 'ats/recruiter/profile/map_update/1' from '192.168.146.46'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 10] [WCN.Role.Session 347] Loading session flash
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Controller.Root 219] ** About to return 1 from root auto
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Controller.ATS 61] ** Enter ATS auto
[DEBUG] [47152f7f] [2012/12/20 15:41:01 1] [WCN.Controller.ATS 145] User 'WCN::DBIC::User::Recruiter=HASH(0x7f9b16810610)' (id: 1) logged in
[DEBUG] [47152f7f] [2012/12/20 15:41:01 1] [WCN.Controller.ATS 950] Validate ATS access rights for recruiter '1' to path 'ats/recruiter/profile/map_update/1'
[WARN] [47152f7f] [2012/12/20 15:41:01 37] [WCN.Controller.ATS 988] User '' does not have access to path 'ats/recruiter/profile/map_update' - ACCESS DENIED
[ERROR] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.Controller.BadRequest 104] 403 FORBIDDEN
[DEBUG] [47152f7f] [2012/12/20 15:41:01 1] [WCN.Controller.Root 324] **** enter root controller's end() method
[INFO] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Controller.Root 339] Have already set status (403) and set a body, so will not render any templates
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Controller.Root 376] Set Content-Length header => '1' bytes
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 231] PERFORMANCE: dispatch finished at 0.20909595489502 secs (took: 0.093096017837524
[DEBUG] [47152f7f] [2012/12/20 15:41:01 4] [WCN.Role.PerformanceLogger 249] PERFORMANCE: finalize finished at 0.213540077209473 secs (took: 0.00444412231445
profile/map_update/1'
[DEBUG] [47152f7f] [2012/12/20 15:41:01 0] [WCN.Role.PerformanceLogger 266] PERFORMANCE: SQL query count: 'SELECT' => 15, 'SET' => 4
[DEBUG] [47152f7f] [2012/12/20 15:41:01 2] [WCN.Role.PerformanceLogger 274] PERFORMANCE: SAN calls: '0' total: '0' avg: '0'
```


Applications: Finance



Applications: Finance

- Event processing
 - Company goes public
 - Stock drops sharply
- Continuous querying
 - Track stocks with gains of 10% in a day
 - Create alerts for major buy/sell transactions
- Real-time response
 - BUY BUY BUY
 - SELL SELL SELL

Applications: Astronomy



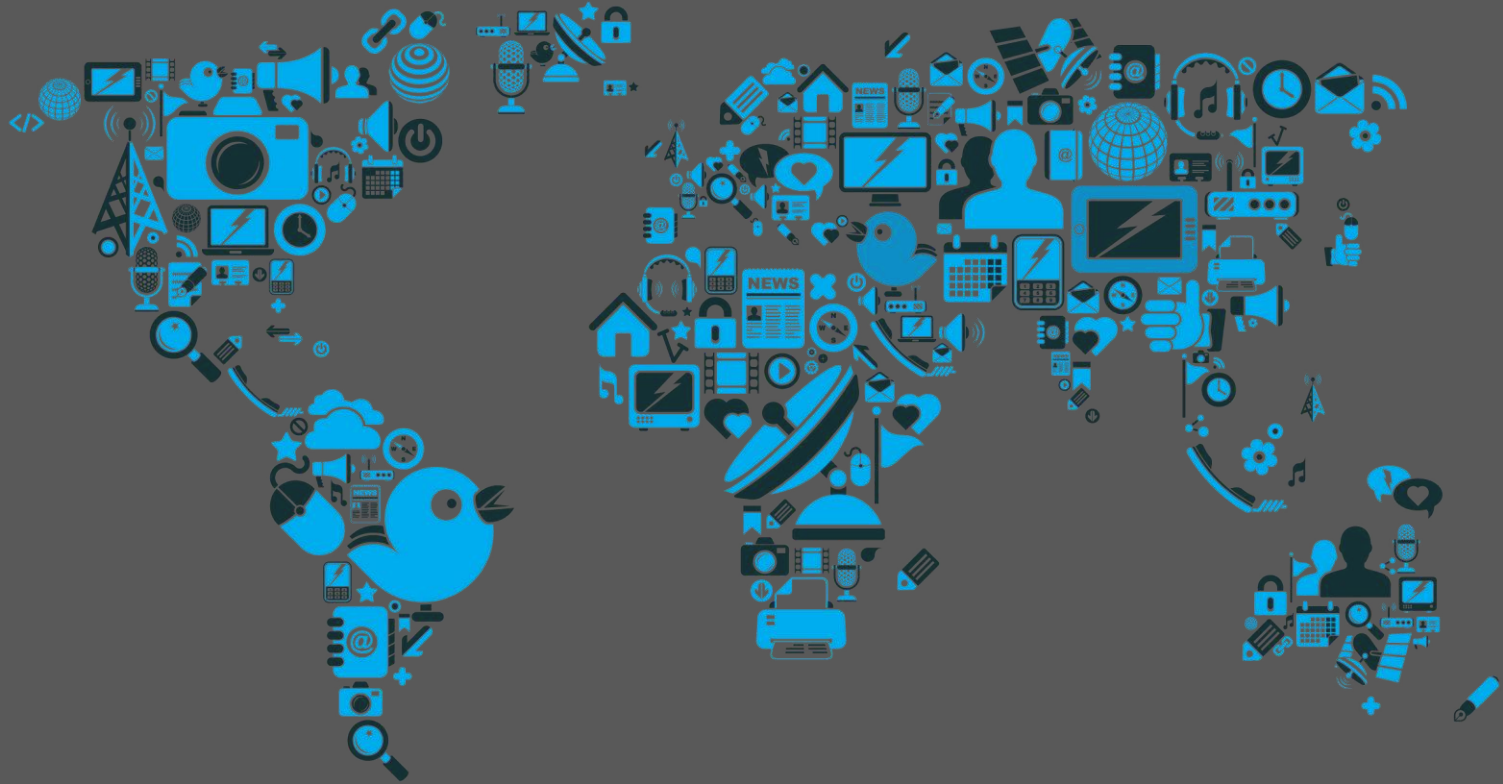
Applications: Astronomy

- Event processing
 - The telescope moves
 - A light source flashes
- Continuous querying
 - Find possible supernovae
 - Track object across the sky
- Real-time response
 - Refocus telescope on important object
 - Lower data filter thresholds

Applications: Astronomy

- Event processing
 - The telescope moves
 - A light source flashes
- Continuous querying
 - Find possible supernovae
 - Track object across the sky
- Real-time response
 - Refocus telescope on important object
 - Lower data filter thresholds

Streams: Internet of Things



Streams: Internet of Things

- Event processing
 - A light turns on
 - It starts to rain
- Continuous querying
 - Tell me when temperature reaches 30°
 - Update position of vehicle
- Real-time response
 - Turn off air conditioning
 - Take another route

DISTRIBUTED STREAMING PLATFORM

Available Frameworks



Application: Emergency Response



chile natural disasters



All **Images** Videos News Maps More

Settings Tools

View saved SafeSearch

tsunami

natural hazard

earthquake

volcano

landslide

mudslide

school

photography

building

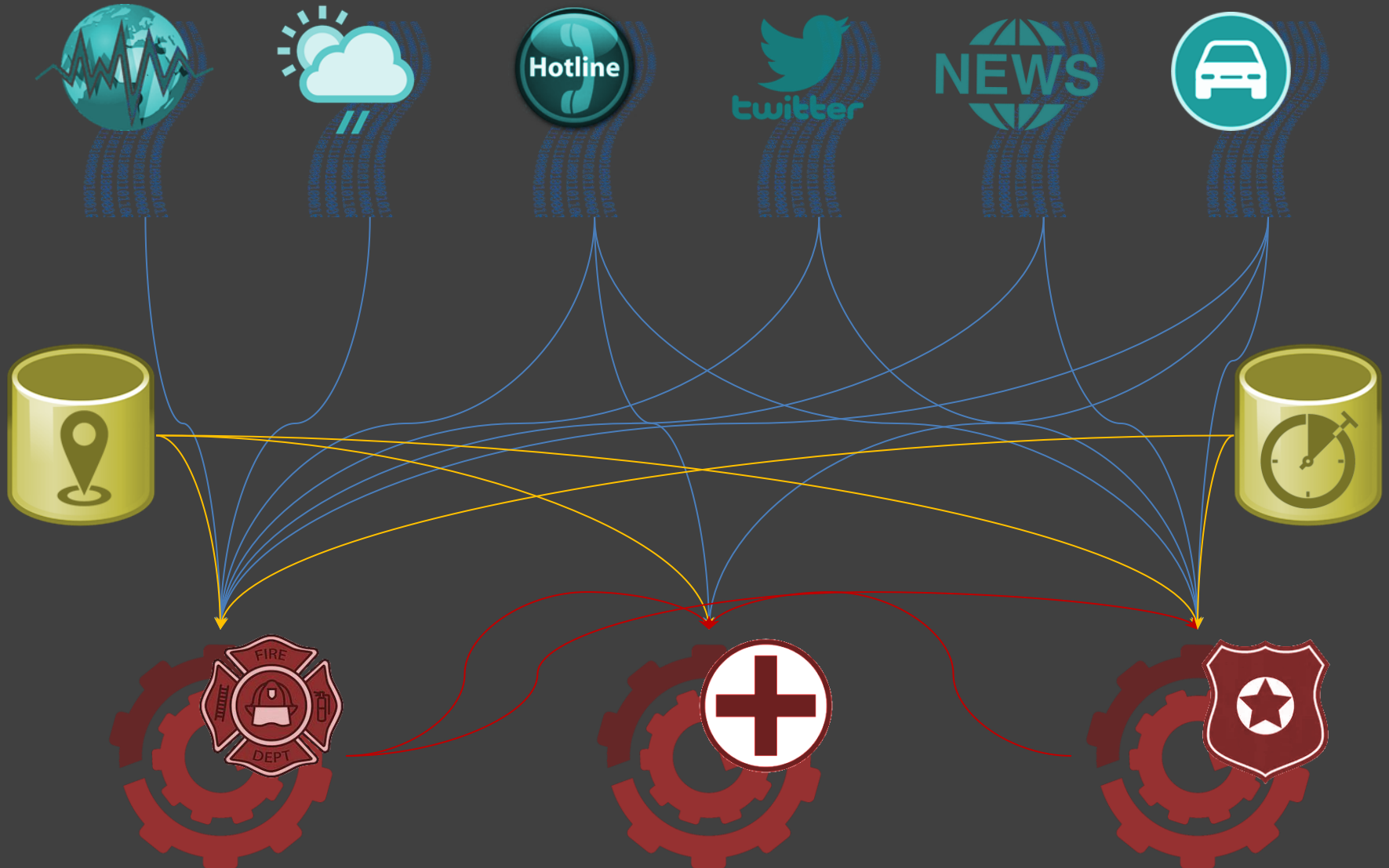
population

historic

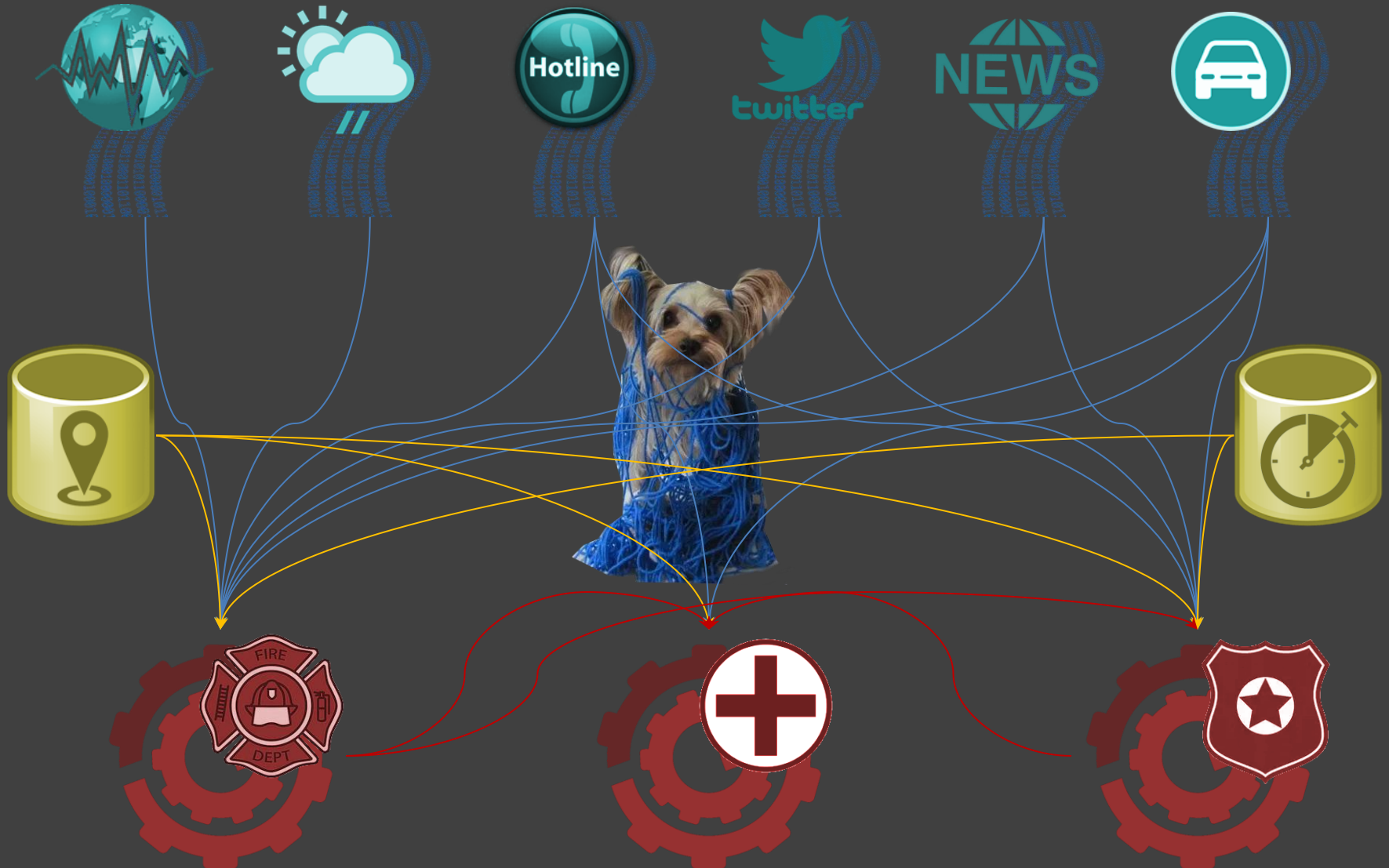
feature



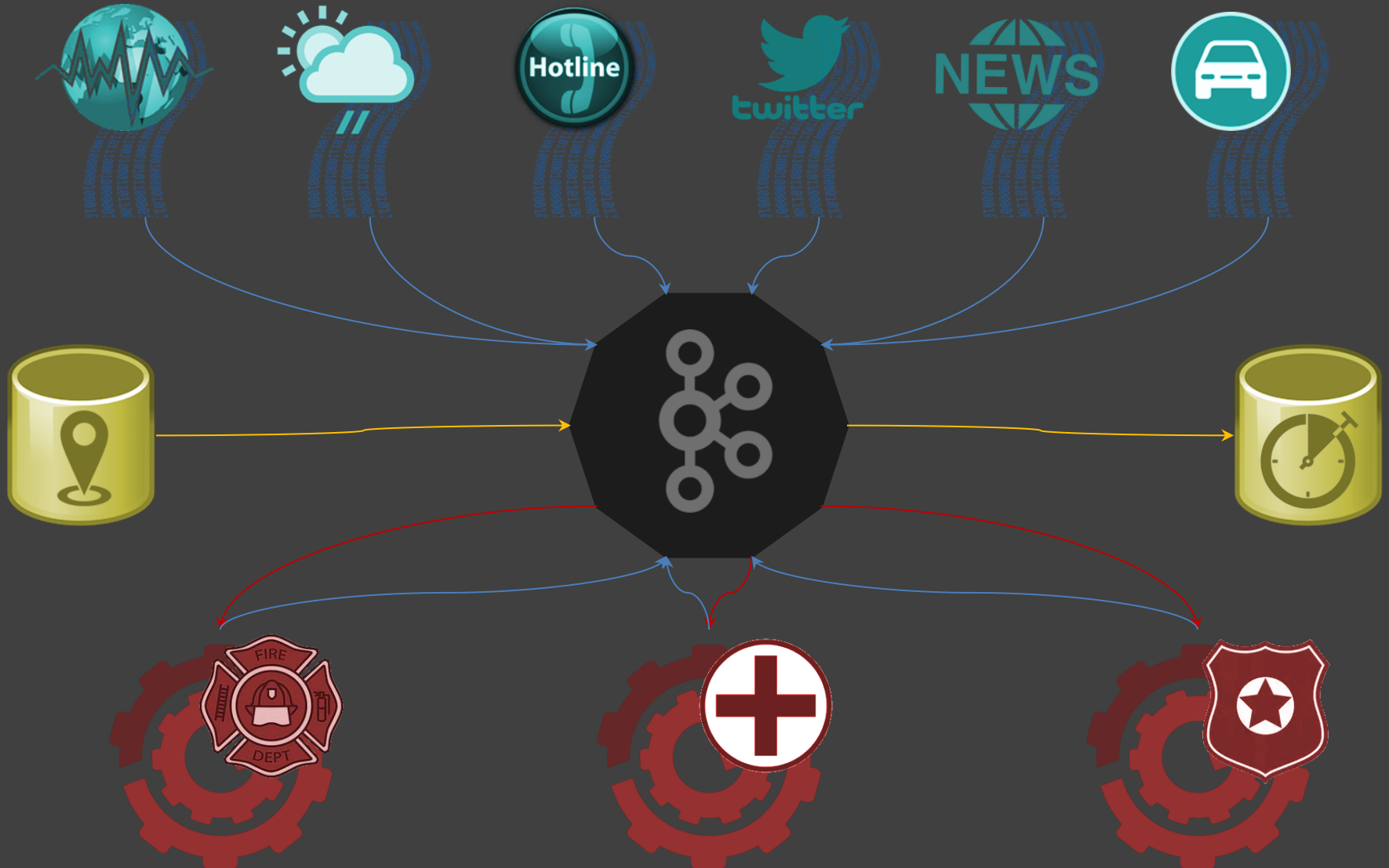
Real-Time Emergency Response



Real-Time Emergency Response



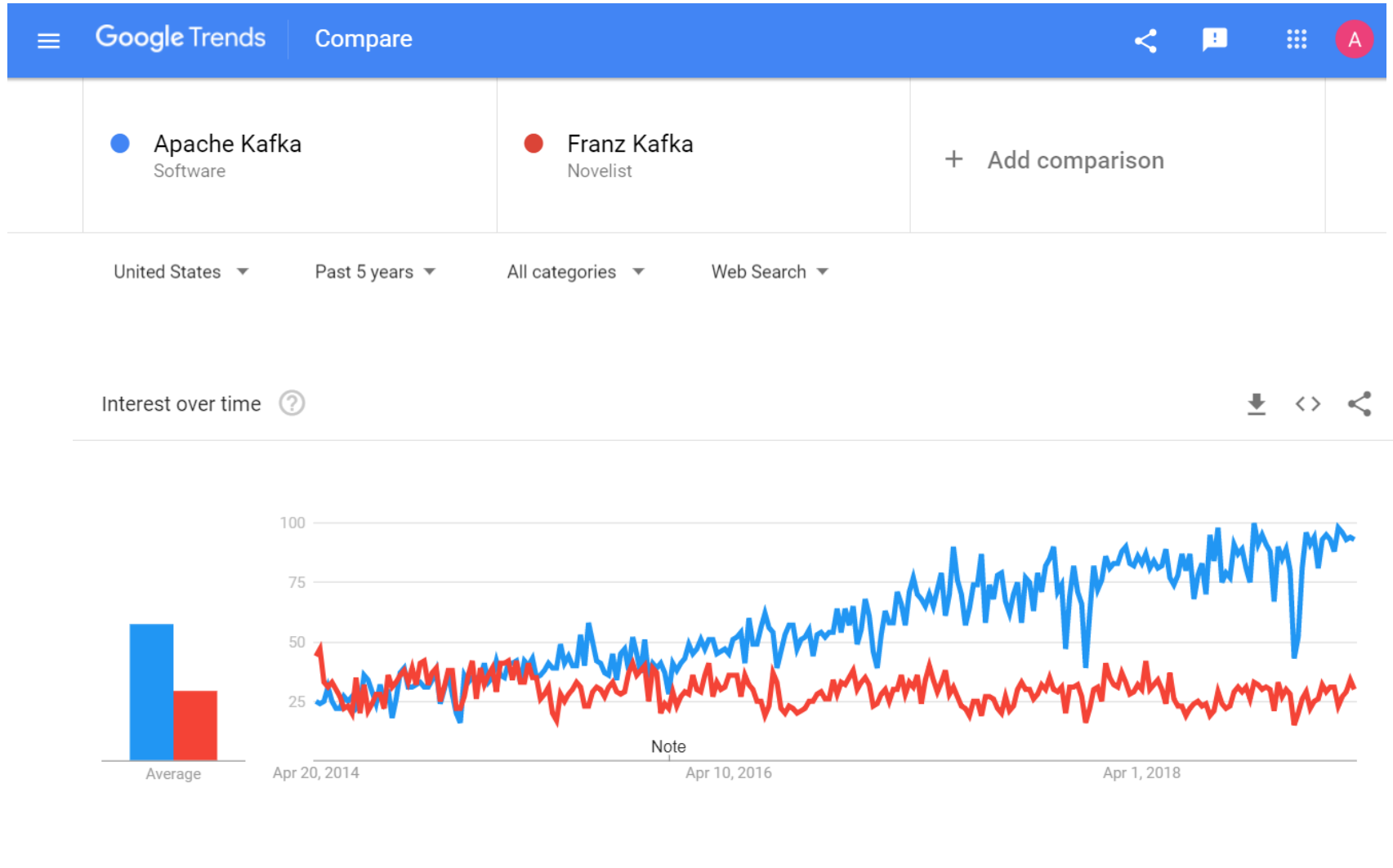
Real-Time Emergency Response



APACHE KAFKA




Apache Kafka vs. Franz Kafka

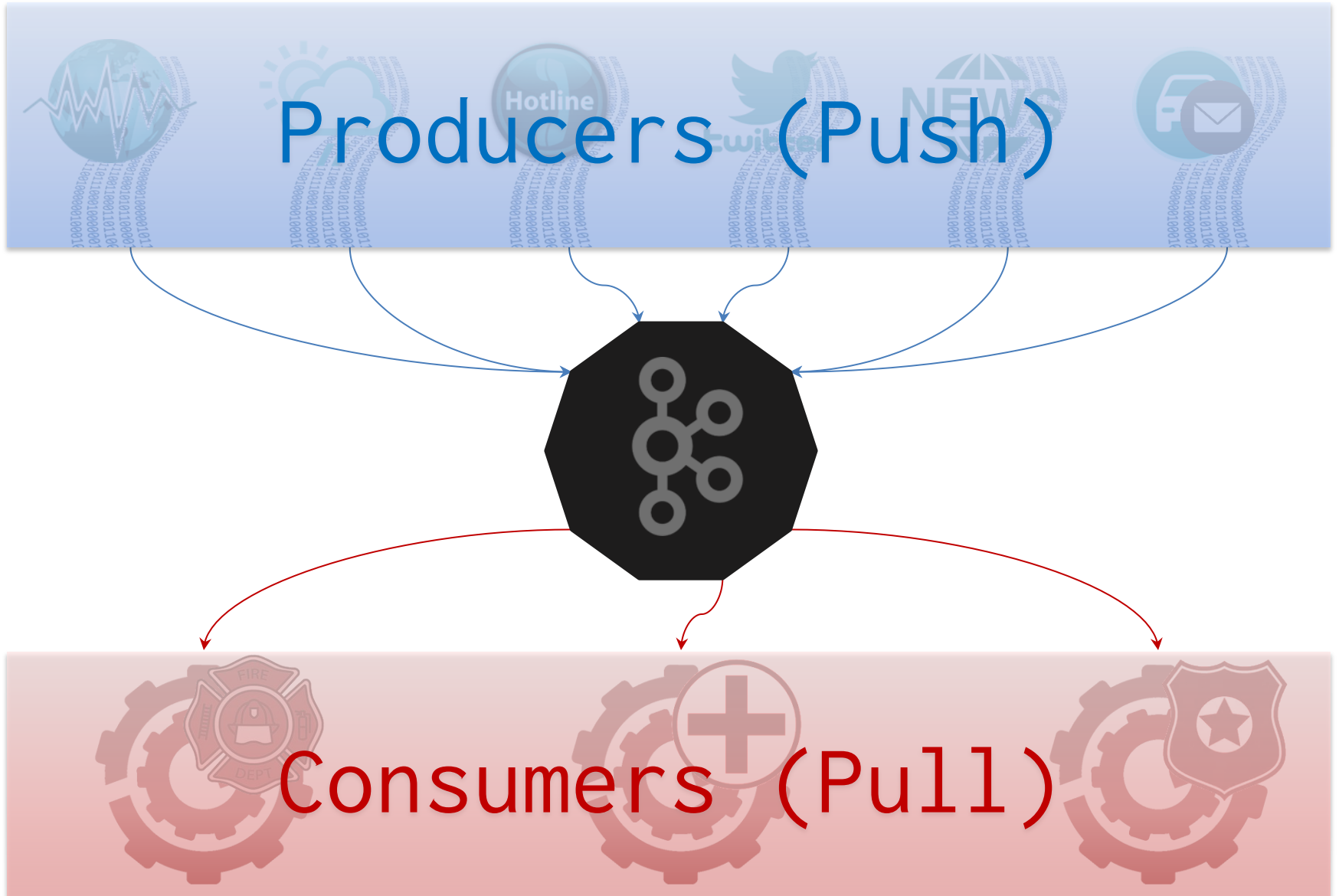


Apache Kafka



- Open Source
- Scala / Java
- Originated in 

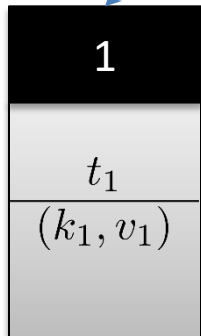
Kafka Overview



KAFKA: DATA MODEL

Kafka Record

Producers



Consumers

Kafka Record

Producers

- Records represent "events"




- Records are immutable

- Contain id (offset), timestamp, key and value
 - Timestamp assigned by application or Kafka

Consumers

Kafka Ledger

Producers



1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

Consumers

Kafka Ledger

Producers

- Producers may only append to ledger



Consumers

Kafka Ledger

Producers

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

Consumers

Kafka Ledger

Producers

- Producers may only append to ledger

1

2

3

4

5

6

7

8

...

t_1

t_2

t_3

t_4

t_5

t_6

t_7

t_8

- Consumers can read from anywhere*

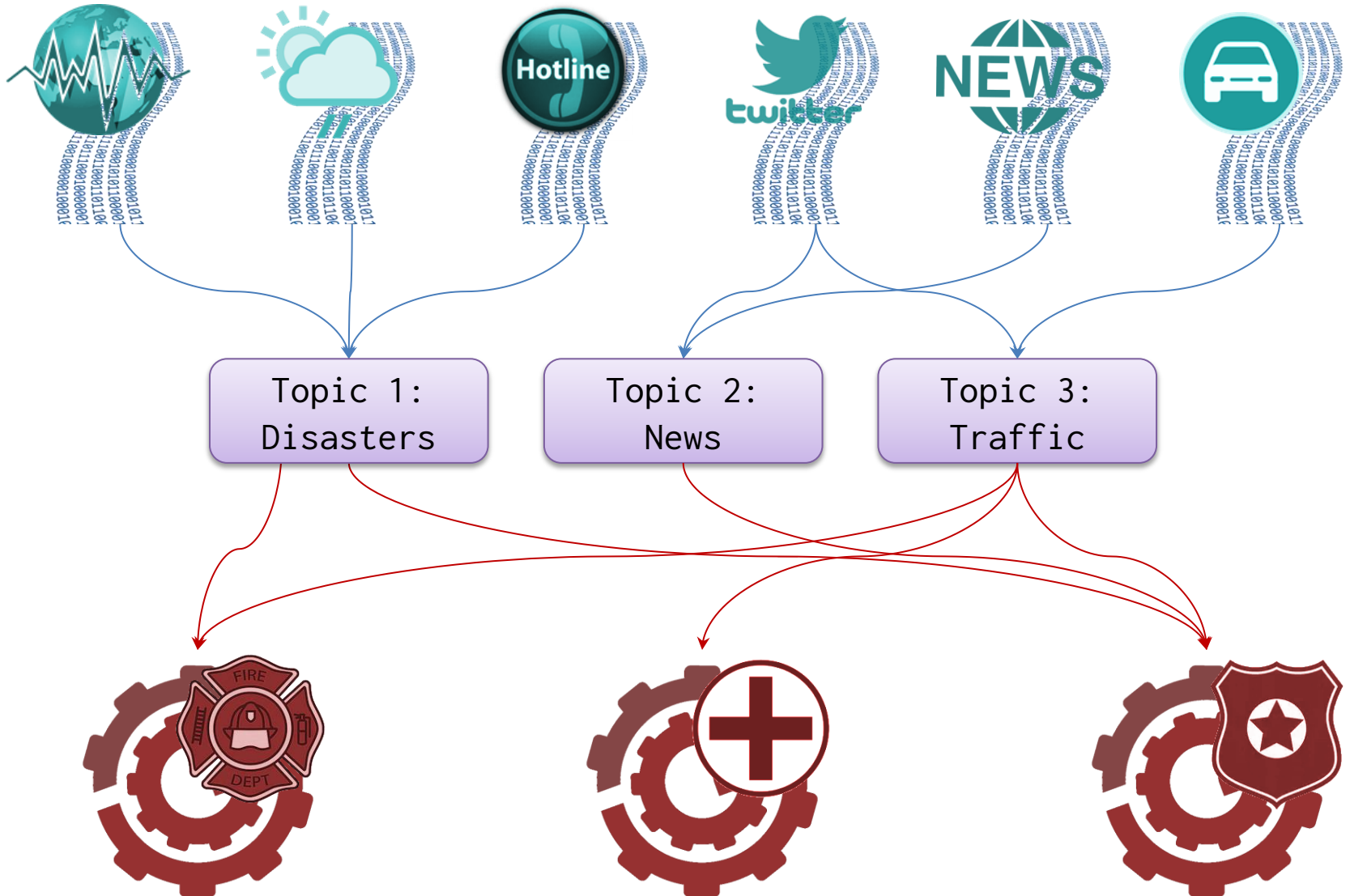
* kind of

Consumers

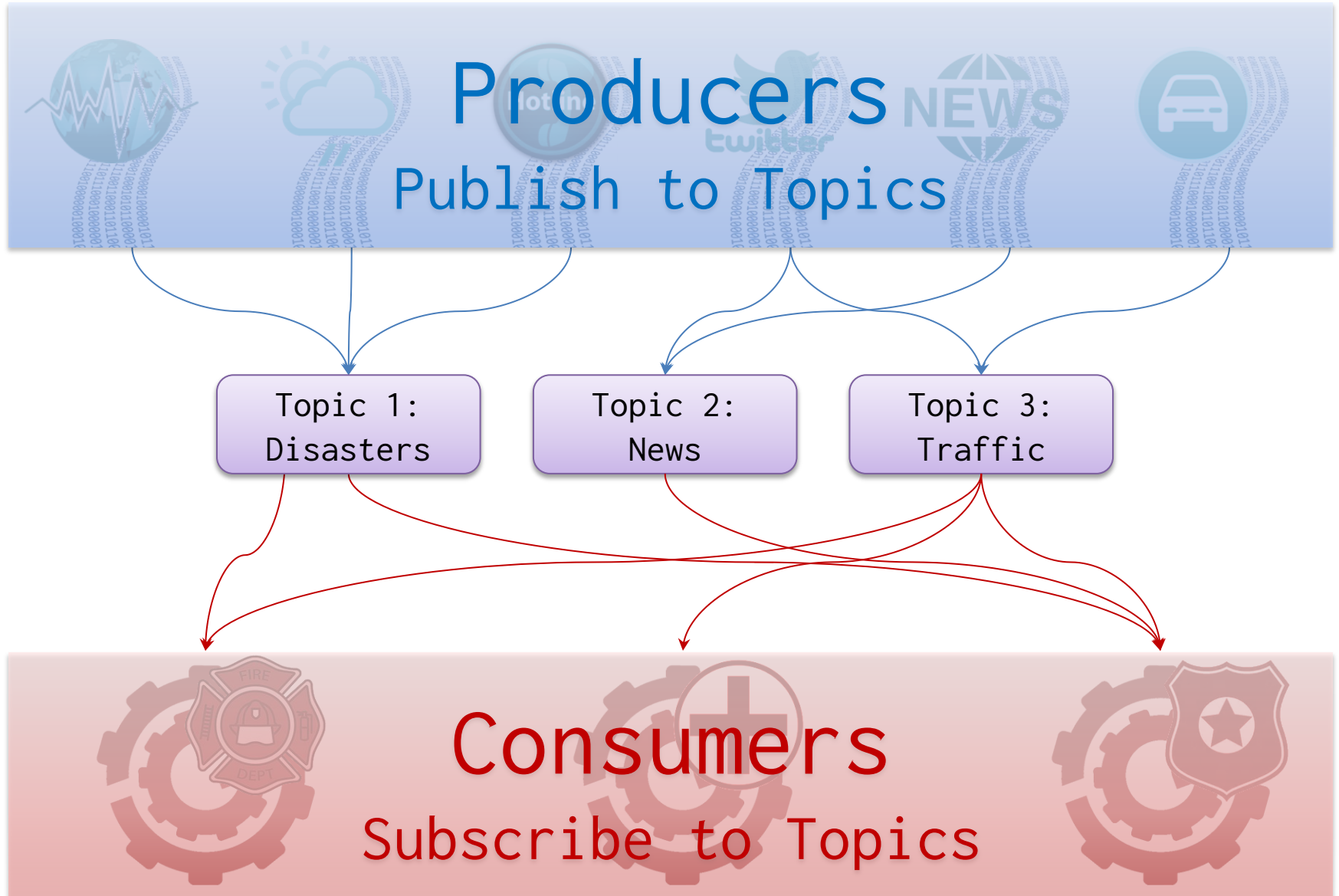


KAFKA: TOPICS

Kafka Topics



Kafka Topics



Topic

Partition 1

- Topics are persistent (on disk)

t_1^1	– Configurable retention policy				
(k_1^1, v_1^1)	(k_2^1, v_2^1)	(k_3^1, v_3^1)	(k_4^1, v_4^1)	(k_5^1, v_5^1)	(k_6^1, v_6^1)
<ul style="list-style-type: none">• Keep everything					

Partition 2

- Delete once consumed

1	• Keep for a period of time				6
t_1^2	t_2^2	t_3^2	t_4^2	t_5^2	t_6^2
(k_1^2, v_1^2)	(k_2^2, v_2^2)	(k_3^2, v_3^2)	(k_4^2, v_4^2)	(k_5^2, v_5^2)	(k_6^2, v_6^2)
<ul style="list-style-type: none">• Use fixed amount of space					

Partition 3


1	2	3	4
t_1^3	t_2^3	t_3^3	t_4^3
(k_1^3, v_1^3)	(k_2^3, v_2^3)	(k_3^3, v_3^3)	(k_4^3, v_4^3)



Topic: Default Partitioning by Key


Partition 1

1	2	3	4	5	6	...
t_1^1	t_2^1	t_3^1	t_4^1	t_5^1	t_6^1	
(k_1^1, v_1^1)	(k_2^1, v_2^1)	(k_3^1, v_3^1)	(k_4^1, v_4^1)	(k_5^1, v_5^1)	(k_6^1, v_6^1)	




Partition 2

1	2	3	4	5	6	7	8	...
t_1^2	t_2^2	t_3^2	t_4^2	t_5^2	t_6^2	t_7^2	t_8^2	
(k_1^2, v_1^2)	(k_2^2, v_2^2)	(k_3^2, v_3^2)	(k_4^2, v_4^2)	(k_5^2, v_5^2)	(k_6^2, v_6^2)	(k_7^2, v_7^2)	(k_8^2, v_8^2)	



Partition 3

1	2	3	4	...
t_1^3	t_2^3	t_3^3	t_4^3	
(k_1^3, v_1^3)	(k_2^3, v_2^3)	(k_3^3, v_3^3)	(k_4^3, v_4^3)	



Topic: Default Partitioning by Key

Partition 1

- Ordering (offset) guaranteed per partition

t_1^1					t_6^1
(k_1^1, v_1^1)	(k_2^1, v_2^1)	(k_3^1, v_3^1)	(k_4^1, v_4^1)	(k_5^1, v_5^1)	(k_6^1, v_6^1)

- Not across partitions!
- For ordering across partitions, use timestamp

Partition 2

1	2	3	4	5	6	7	8	...
t_1^2	t_2^2	t_3^2	t_4^2	t_5^2	t_6^2	t_7^2	t_8^2	
(k_1^2, v_1^2)	(k_2^2, v_2^2)	(k_3^2, v_3^2)	(k_4^2, v_4^2)	(k_5^2, v_5^2)	(k_6^2, v_6^2)	(k_7^2, v_7^2)	(k_8^2, v_8^2)	

Partition 3

1	2	3	4	...
t_1^3	t_2^3	t_3^3	t_4^3	
(k_1^3, v_1^3)	(k_2^3, v_2^3)	(k_3^3, v_3^3)	(k_4^3, v_4^3)	

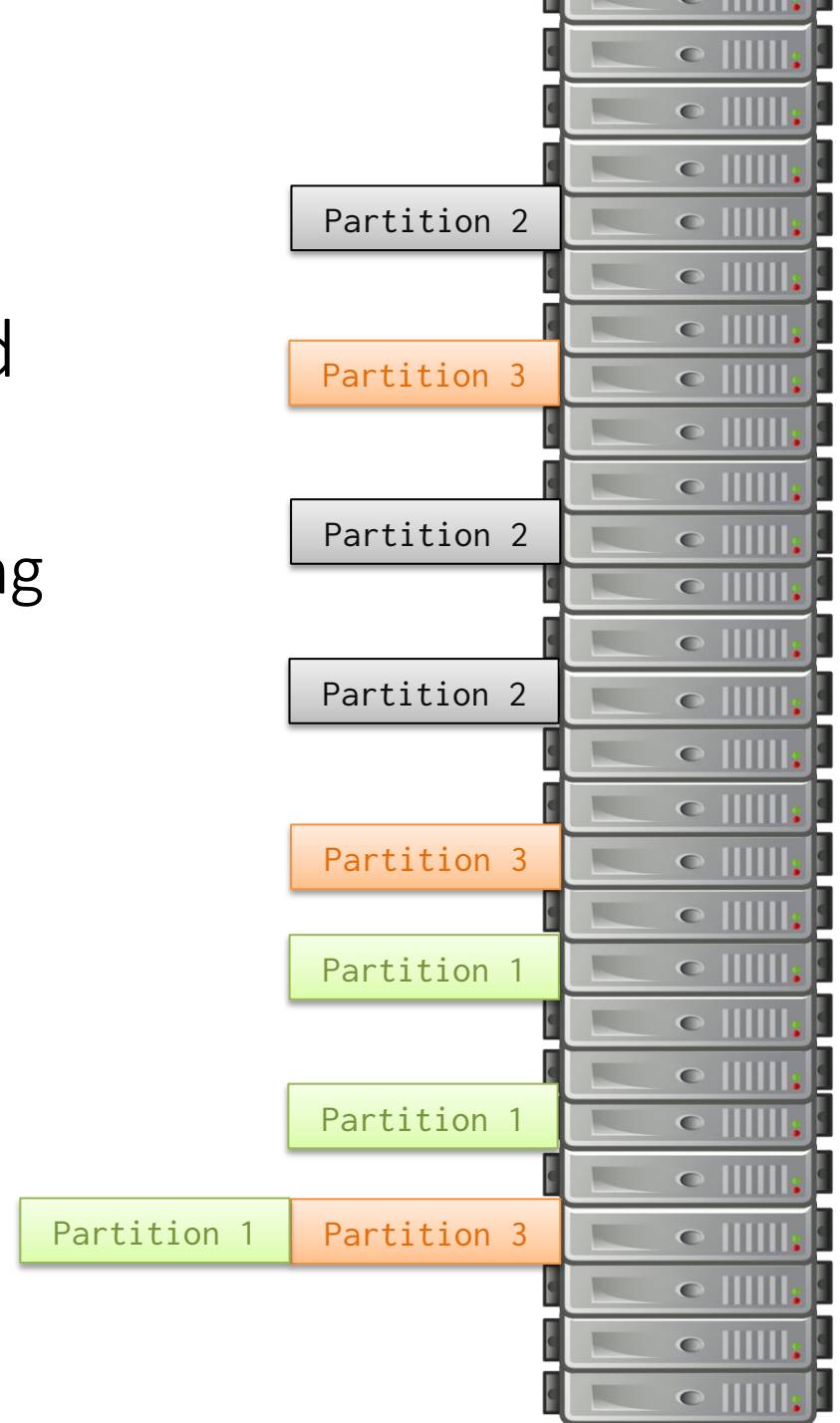
Replication

- Topics can be replicated
 - Choose factor per topic
 - Automatic load balancing

Problem?

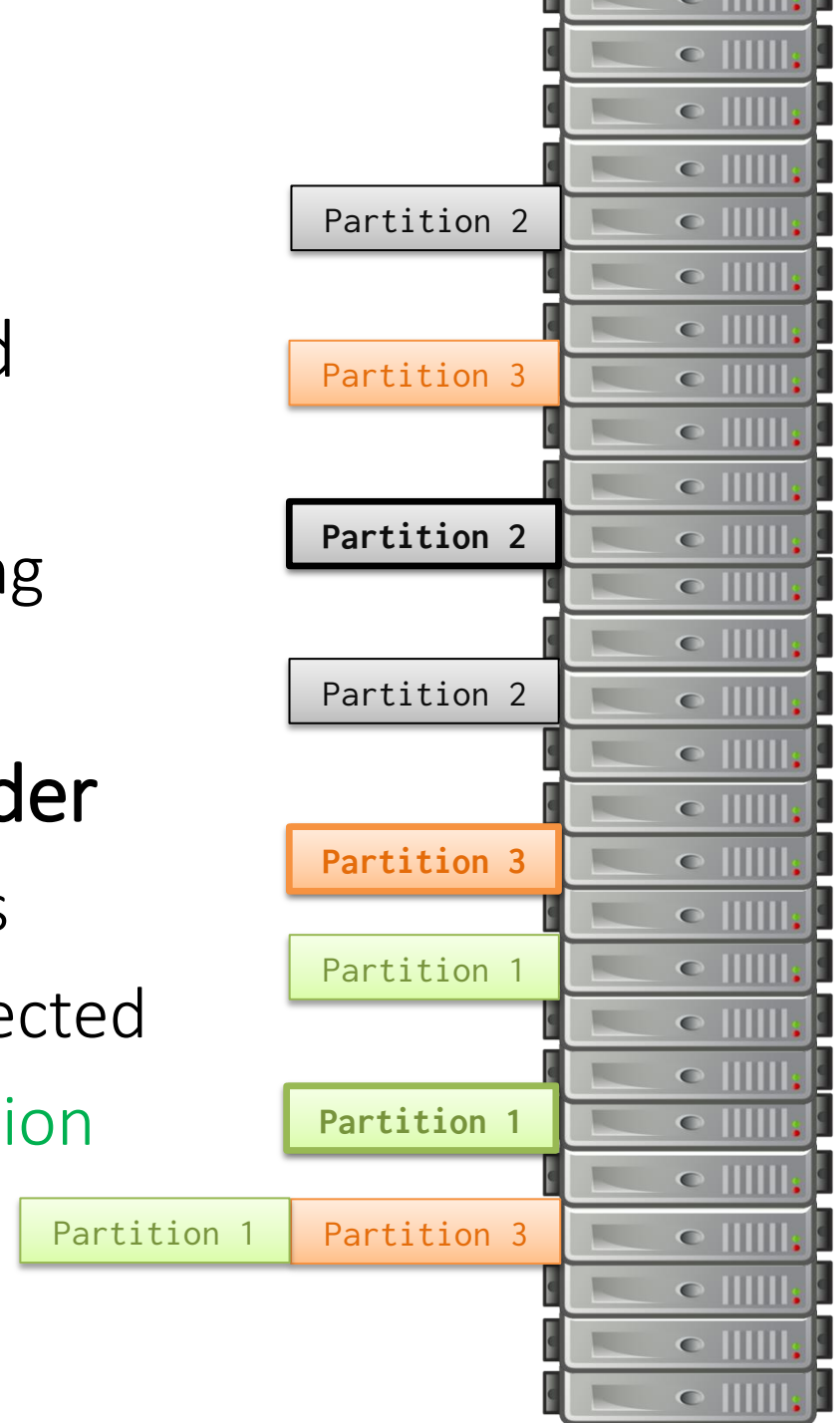


Order?



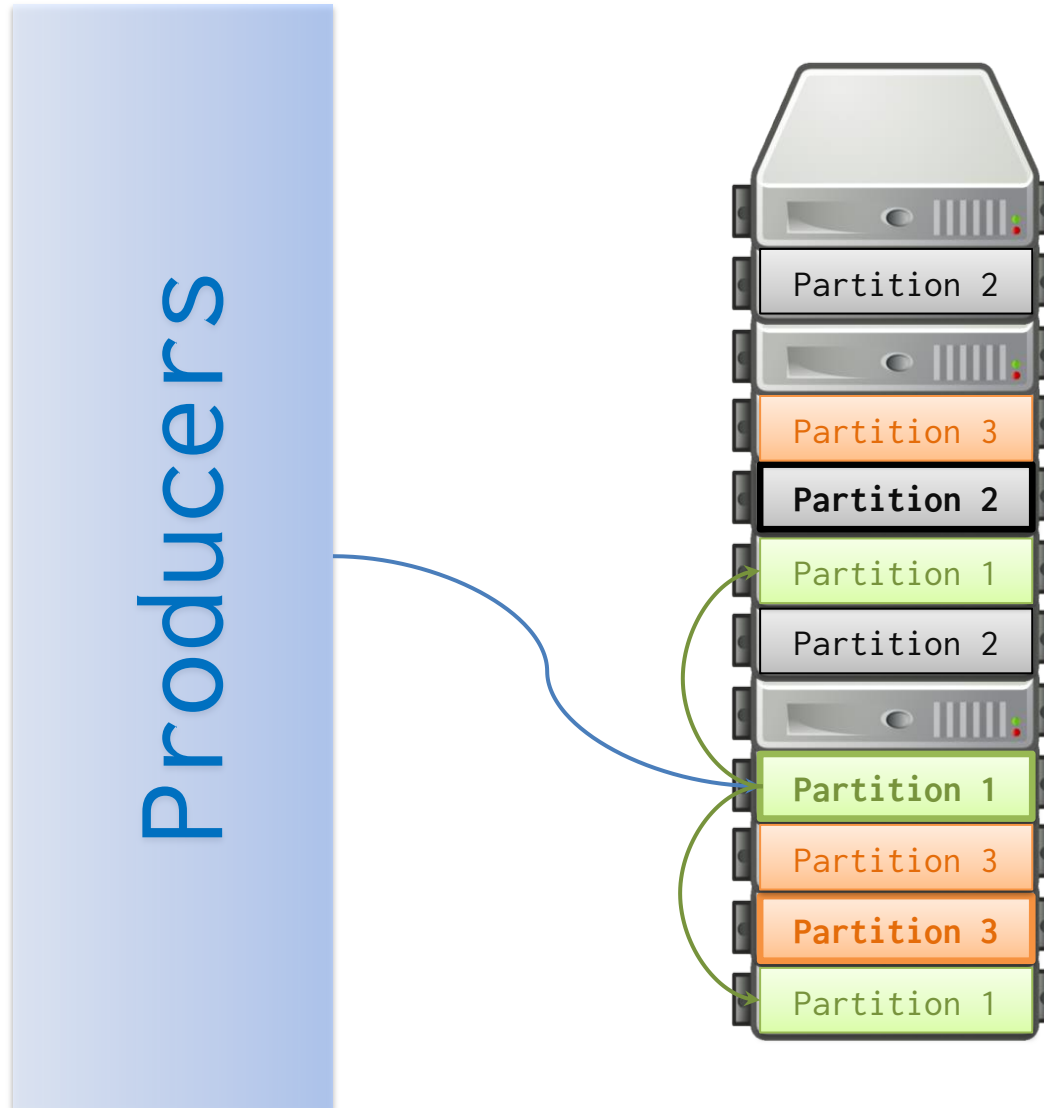
Leader

- Topics can be replicated
 - Choose factor per topic
 - Automatic load balancing
- One machine is the **leader**
 - The others are followers
 - Leader automatically elected
 - Ensures order per partition
 - Reads/writes to leader

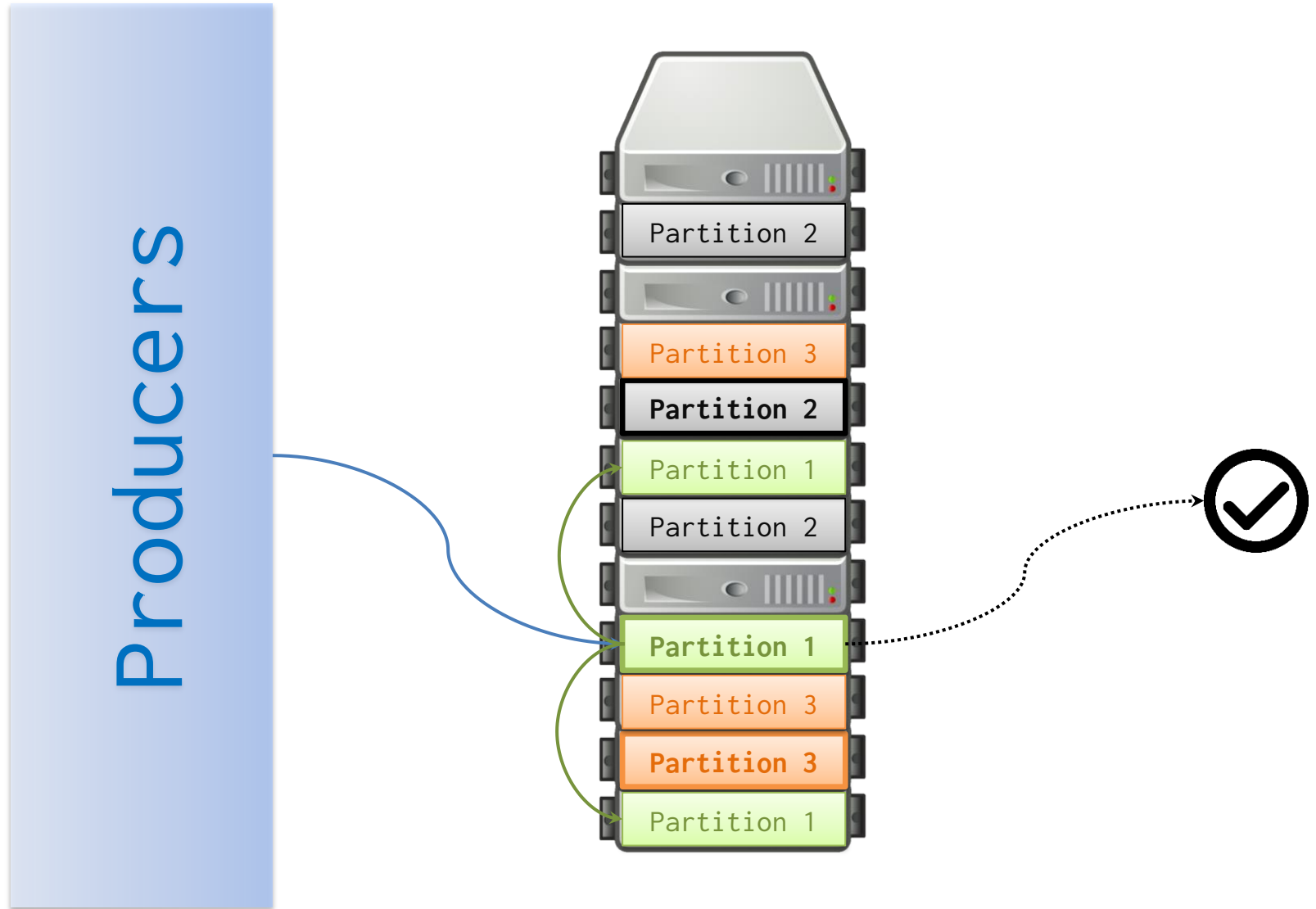


KAFKA: WRITE GUARANTEES

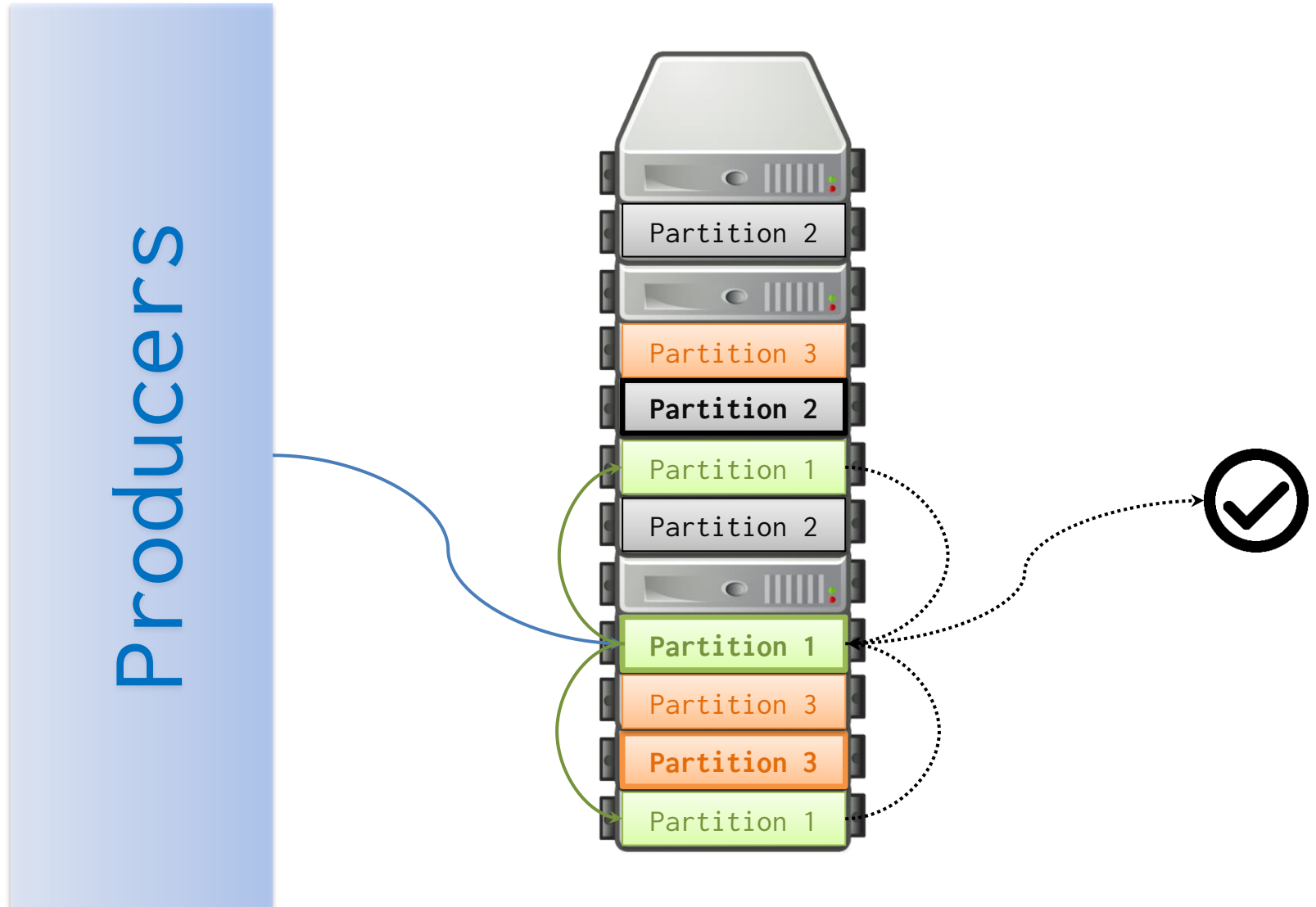
Writes: Asynchronous (No Guarantee)



Writes: Leader Commit



Writes: Leader Commit + Quorum (2)



Write Guarantees

- Asynchronous
 - No guarantee
 - Very low latency
- Leader Commit
 - Persistent on leader
 - Medium latency (disk write + network ack)
- Leader Commit + Quorum n
 - Persistent on leader + n machines
 - High latency (disk writes + network acks)

KAFKA: READS

Kafka tracks consumer offset

C1: 1-2

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



Kafka tracks consumer offset

C1: 3-4

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



1
2
3
4

Kafka tracks consumer offset

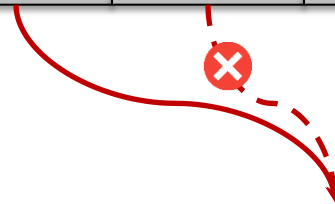
C1: 5-6

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



Failures?

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



- 1
- 2
- 3
- 4
- 5
- 6

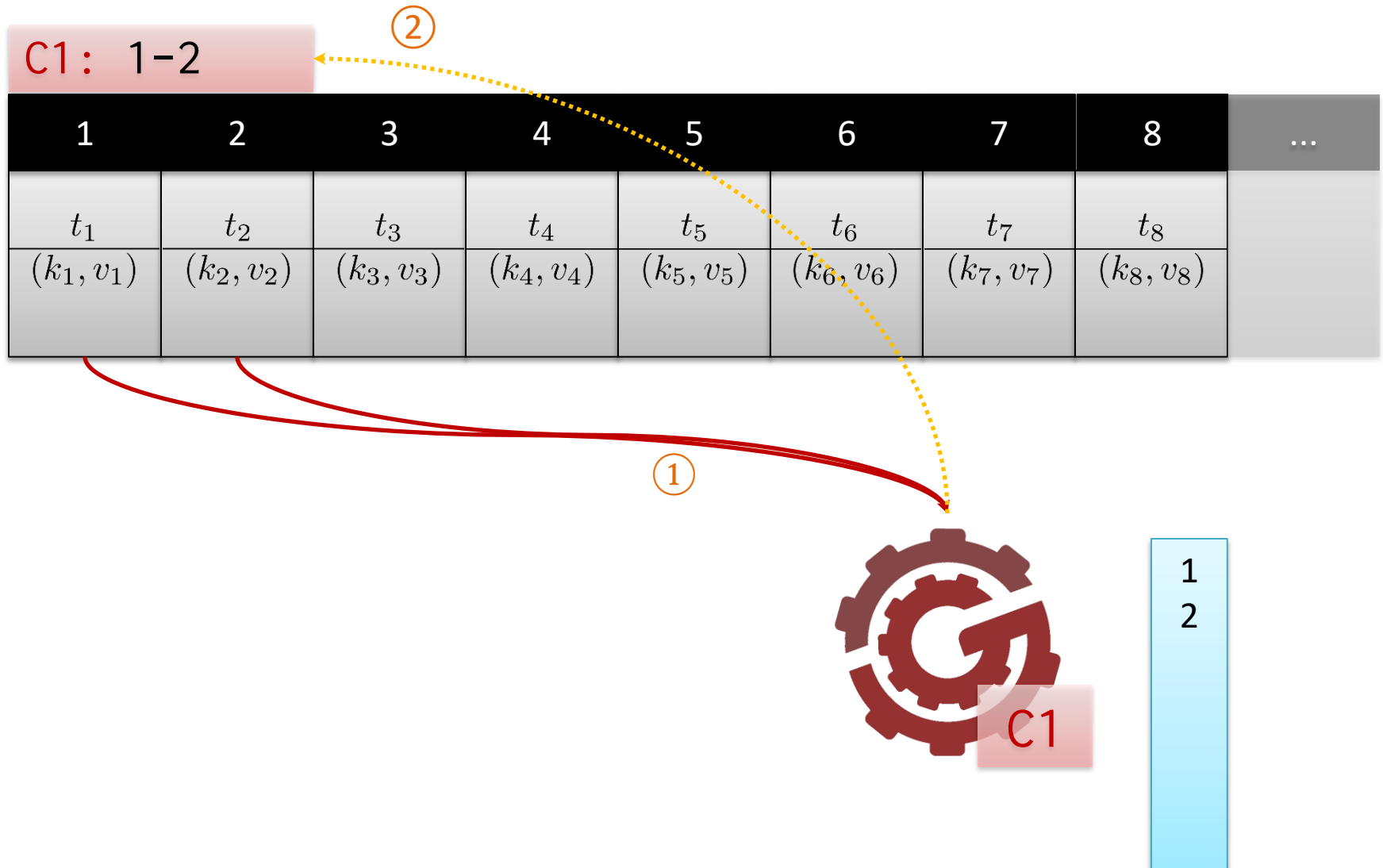
What should we do in the case of a read failure?

KAFKA: READ GUARANTEES

Read Guarantees

- At least once
 - Each value processed at least once
 - Consumer offset updated on consumer ACK
- At most once
- Effectively once
- Exactly once

Read: At Least Once (Default)



Read: At Least Once (Default)

C1: 1-2

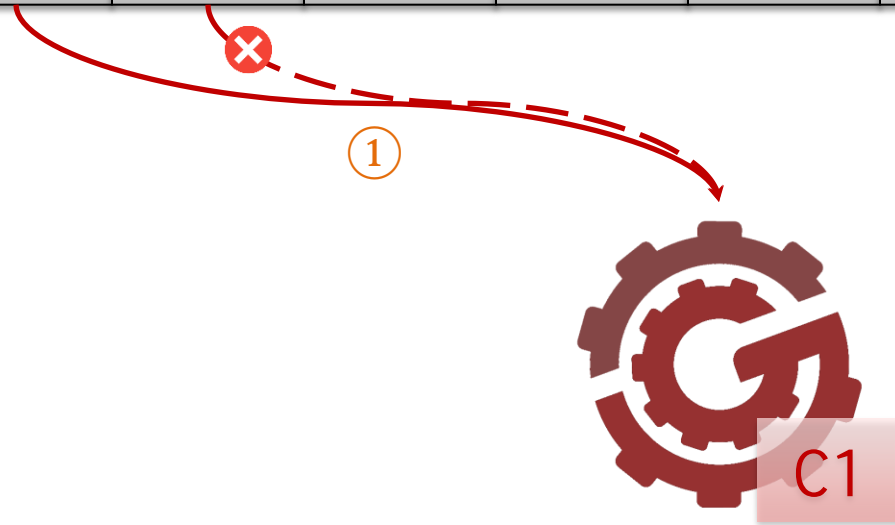
1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



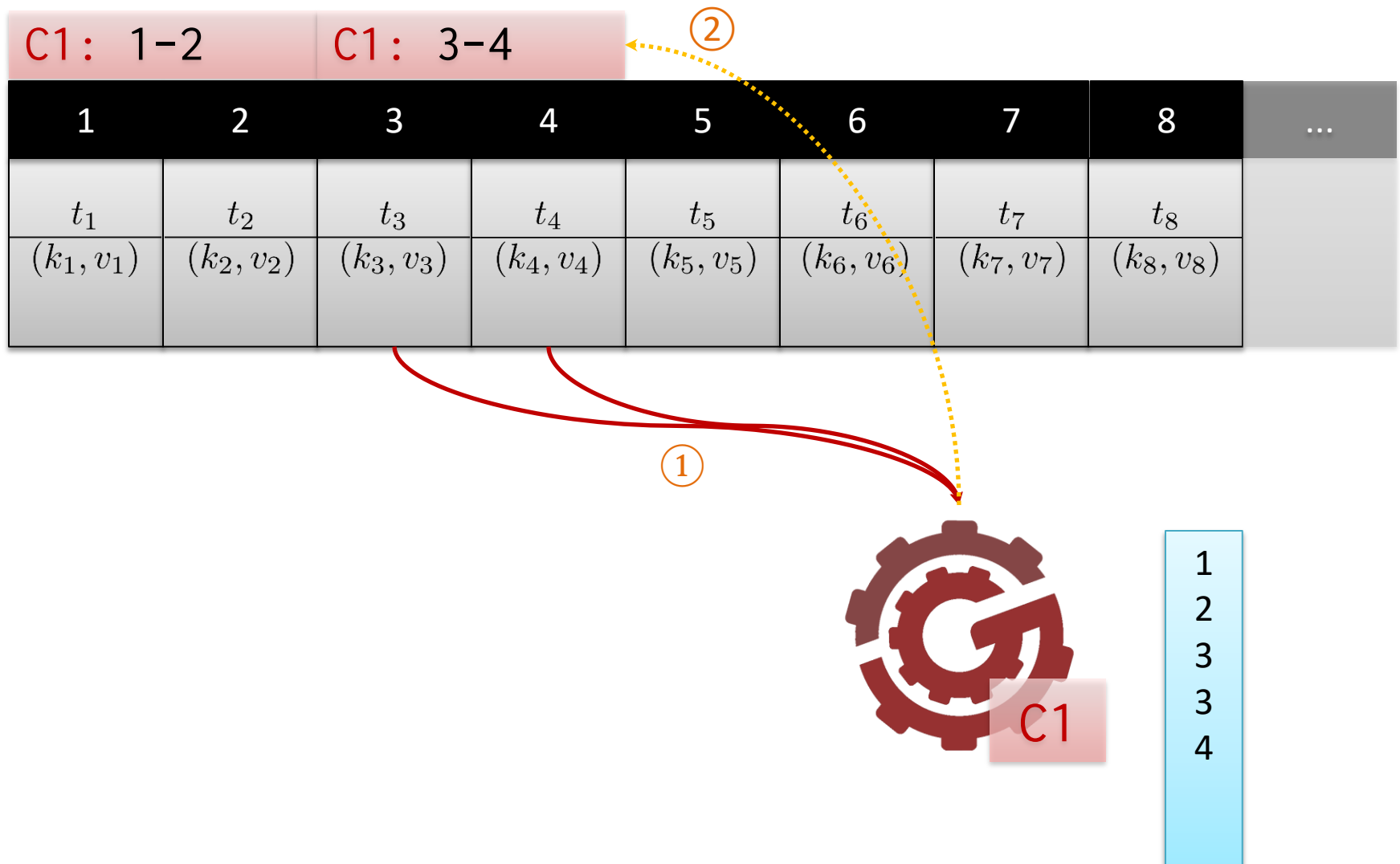
Read: At Least Once (Default)

C1: 1-2

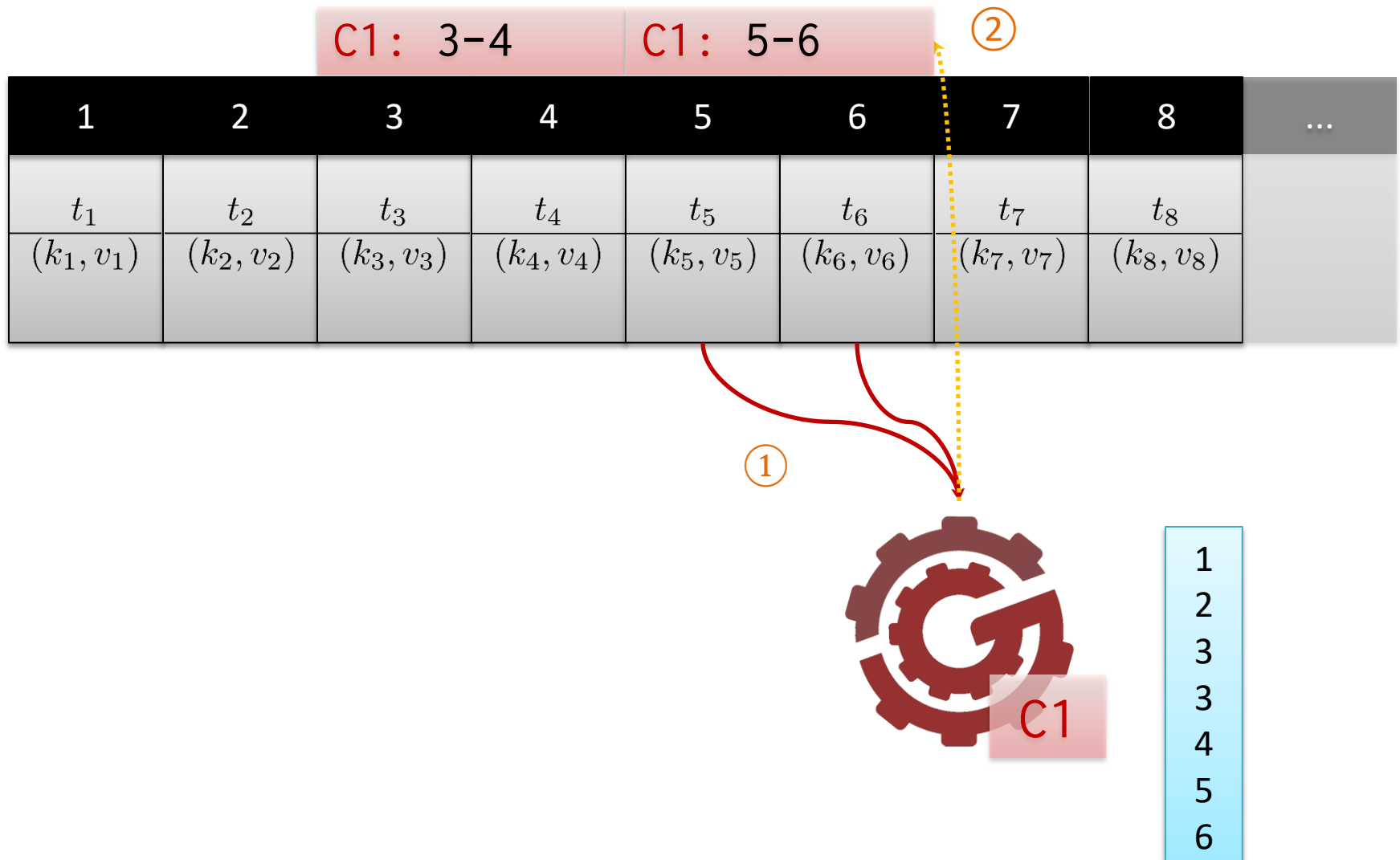
1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



Read: At Least Once (Default)



Read: At Least Once (Default)



Read Guarantees

- At least once
- At most once
 - Each value processed at most once
 - Consumer offset updated immediately
- Effectively once
- Exactly once

Read: At Most Once

①

C1: 1-2

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

②



Read: At Most Once

①

C1: 3-4

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

②



- 1
- 2
- 3

Read: At Most Once

①

C1: 5-6

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

②

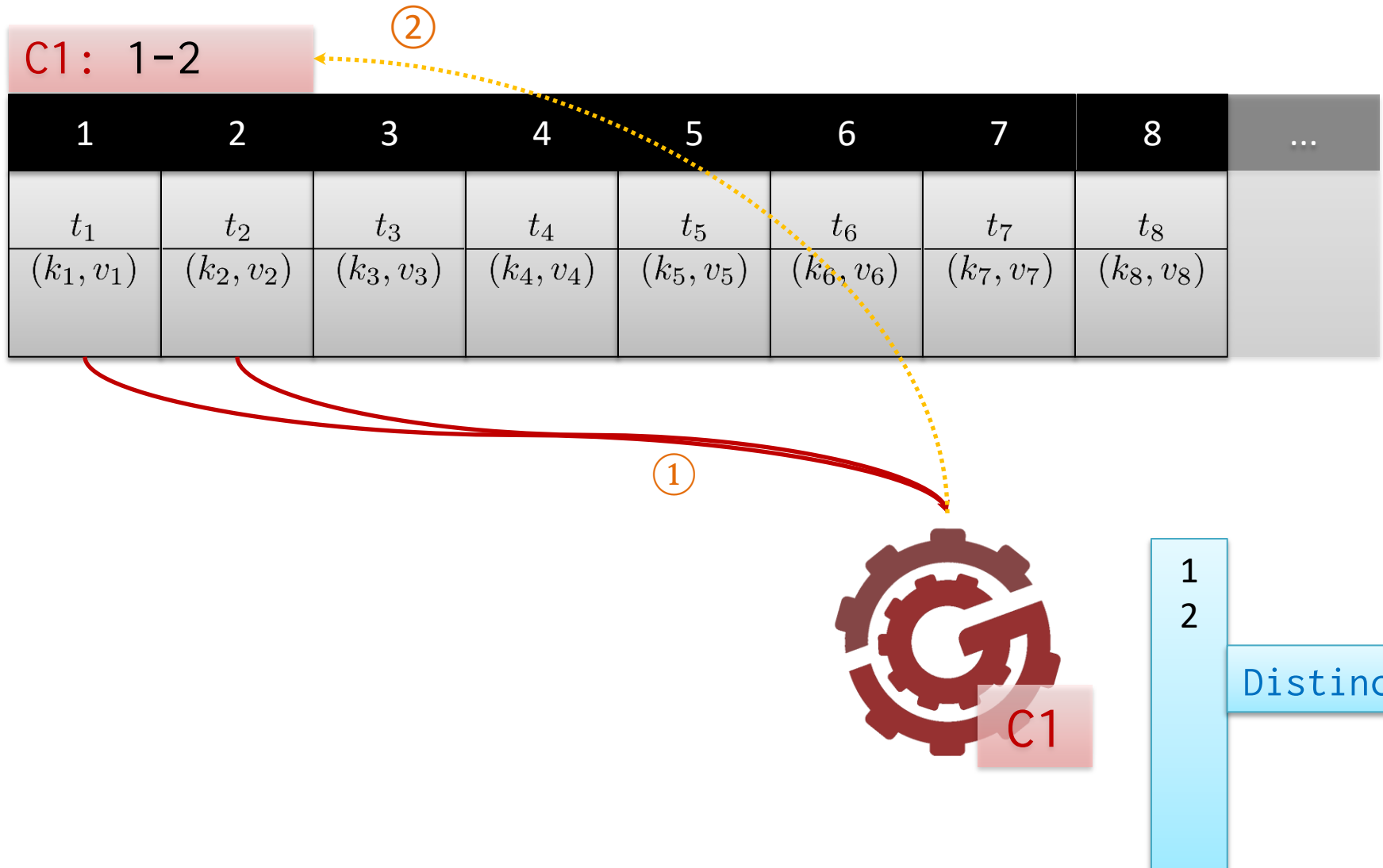


1
2
3
5
6

Read Guarantees

- At least once
- At most once
- Effectively once
 - At least once but ...
 - Consumer takes care of duplicates
- Exactly once

Read: Effectively Once



Read: Effectively Once

C1: 1-2

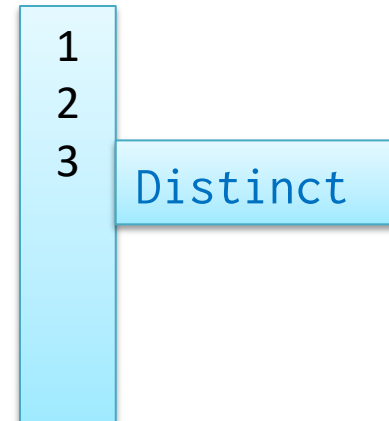
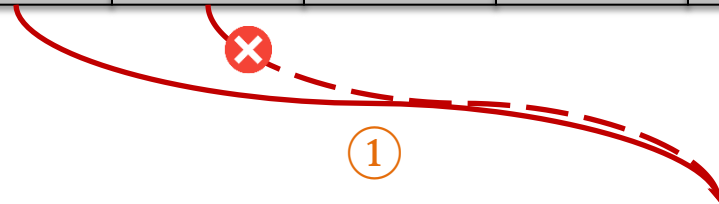
1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



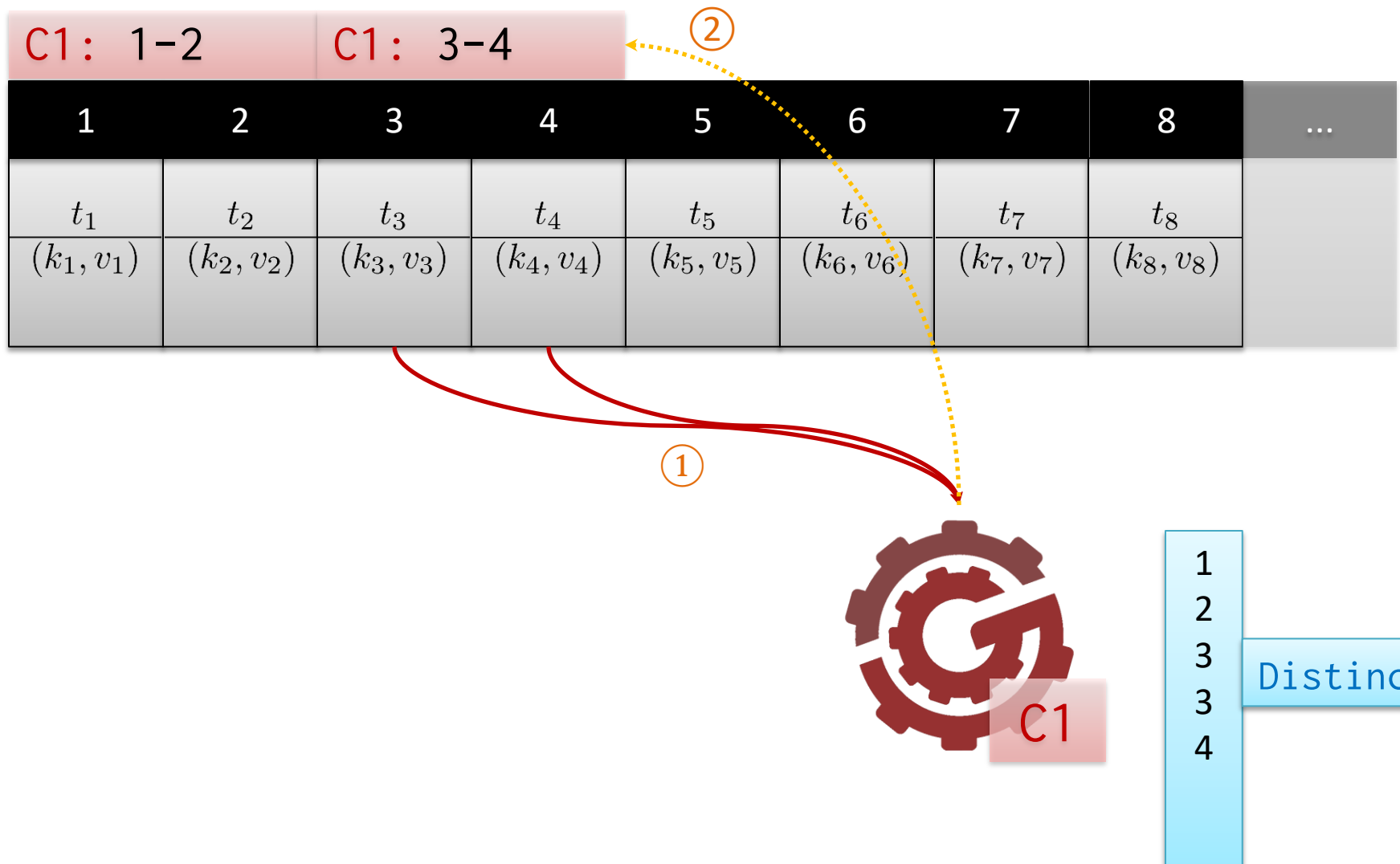
Read: Effectively Once

C1: 1-2

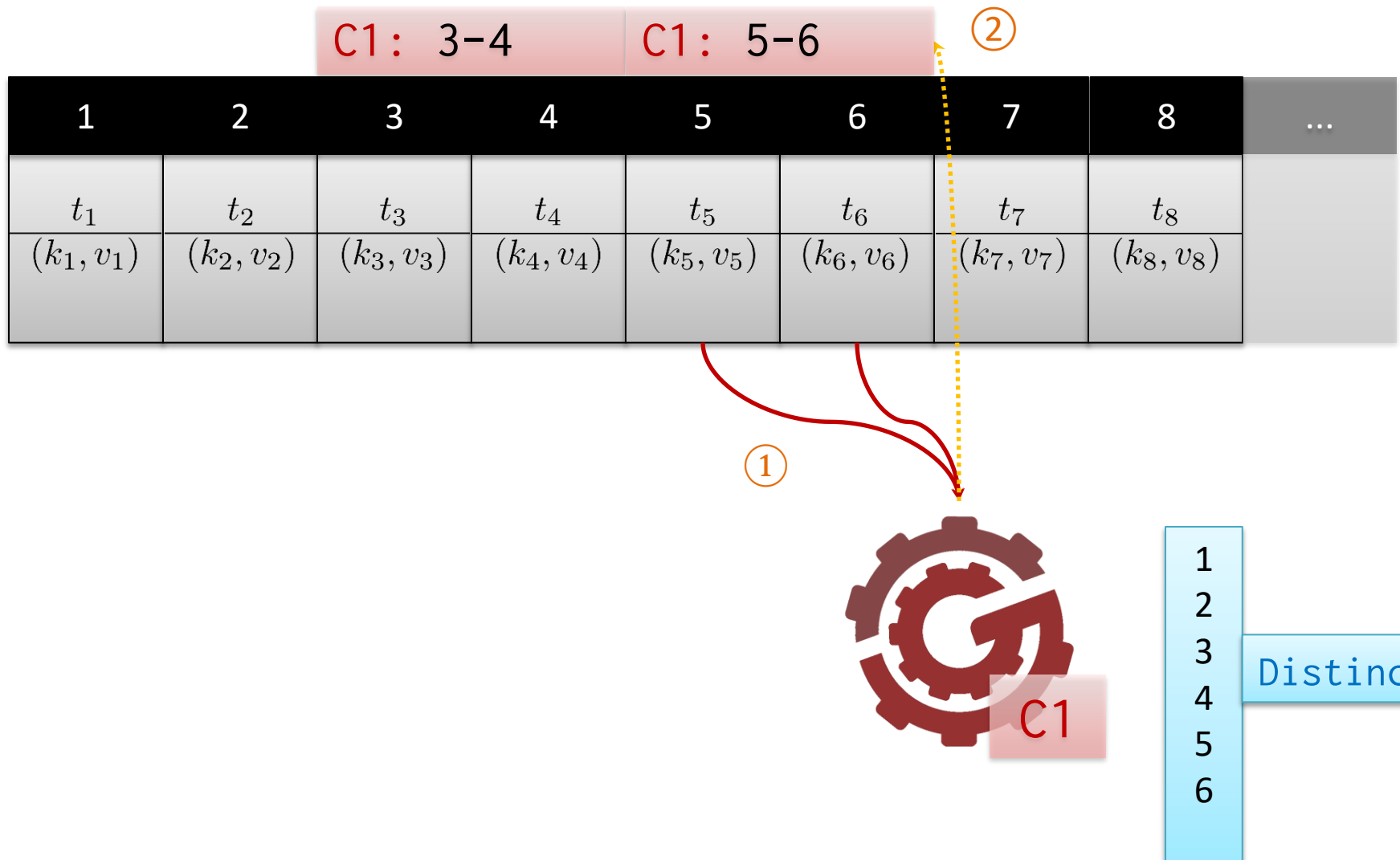
1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



Read: Effectively Once



Read: Effectively Once

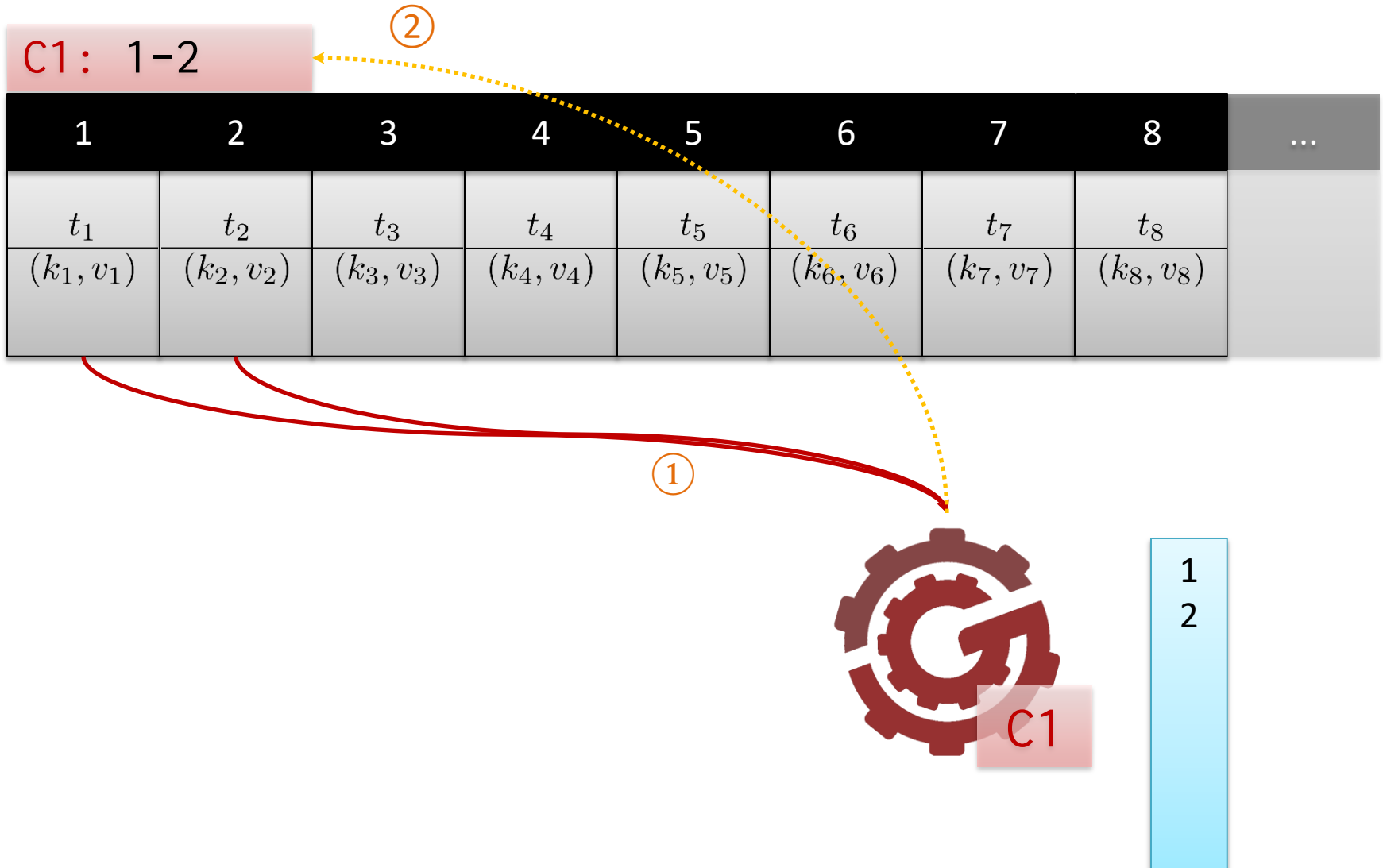


Read Guarantees

- At least once
- At most once
- Effectively once
- Exactly once
 - Data and offset updated as a single transaction

Read: Exactly Once

Transaction 1: ①, ②



Read: Exactly Once

C1: 1-2

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

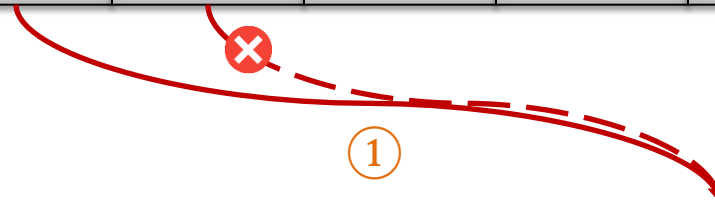


Read: Exactly Once

Transaction 2: ①, ②

C1: 1-2

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	



Read: Exactly Once

Transaction 3: ①, ②

C1: 1-2		C1: 3-4						
1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

①

②



1
2
3
4

Read: Exactly Once

Transaction 4: ①, ②

1	2	3	4	5	6	7	8	...
t_1	t_2	t_3	t_4	t_5	t_6	t_7	t_8	
(k_1, v_1)	(k_2, v_2)	(k_3, v_3)	(k_4, v_4)	(k_5, v_5)	(k_6, v_6)	(k_7, v_7)	(k_8, v_8)	

C1: 3-4

C1: 5-6

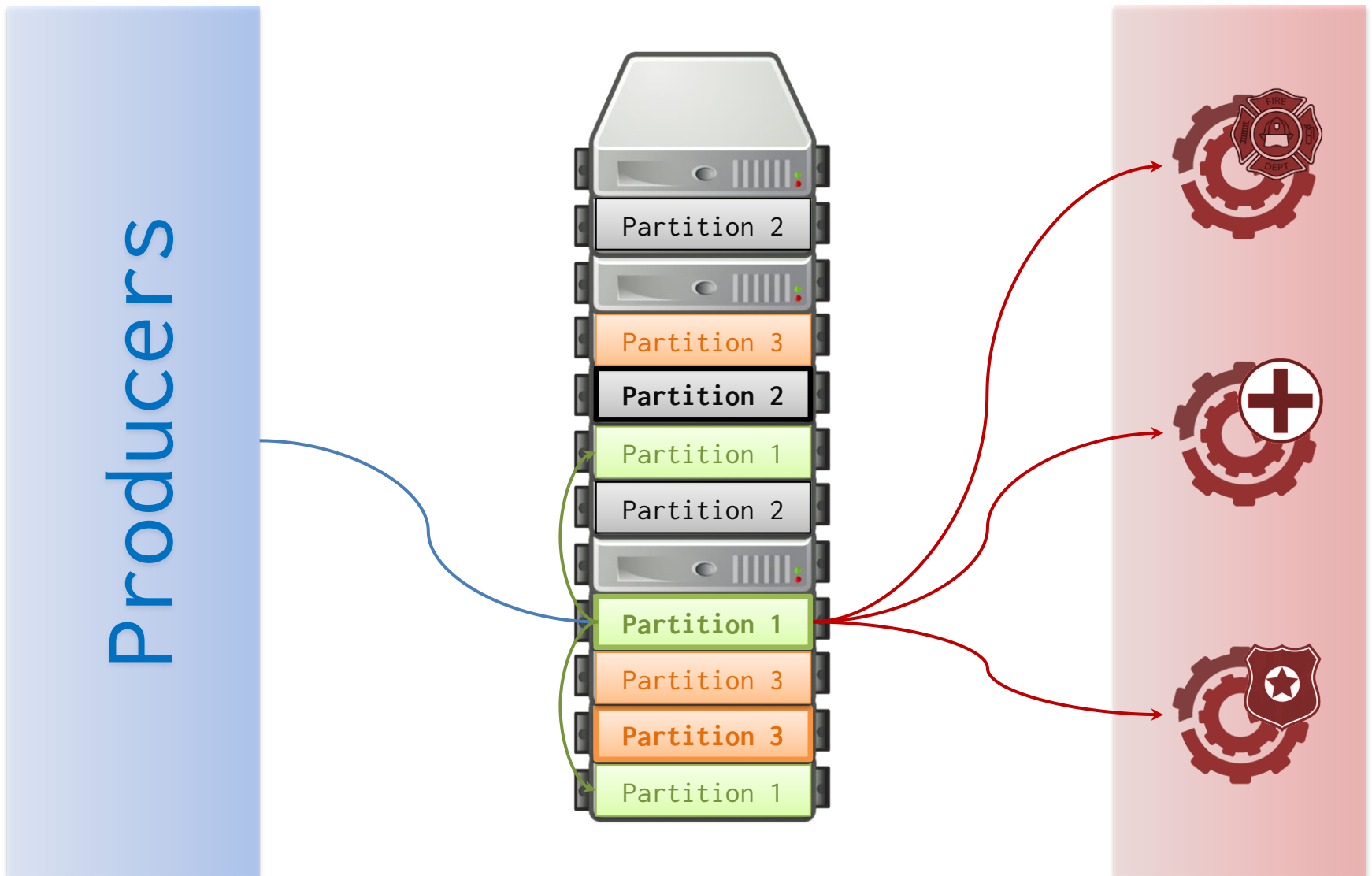
②

①



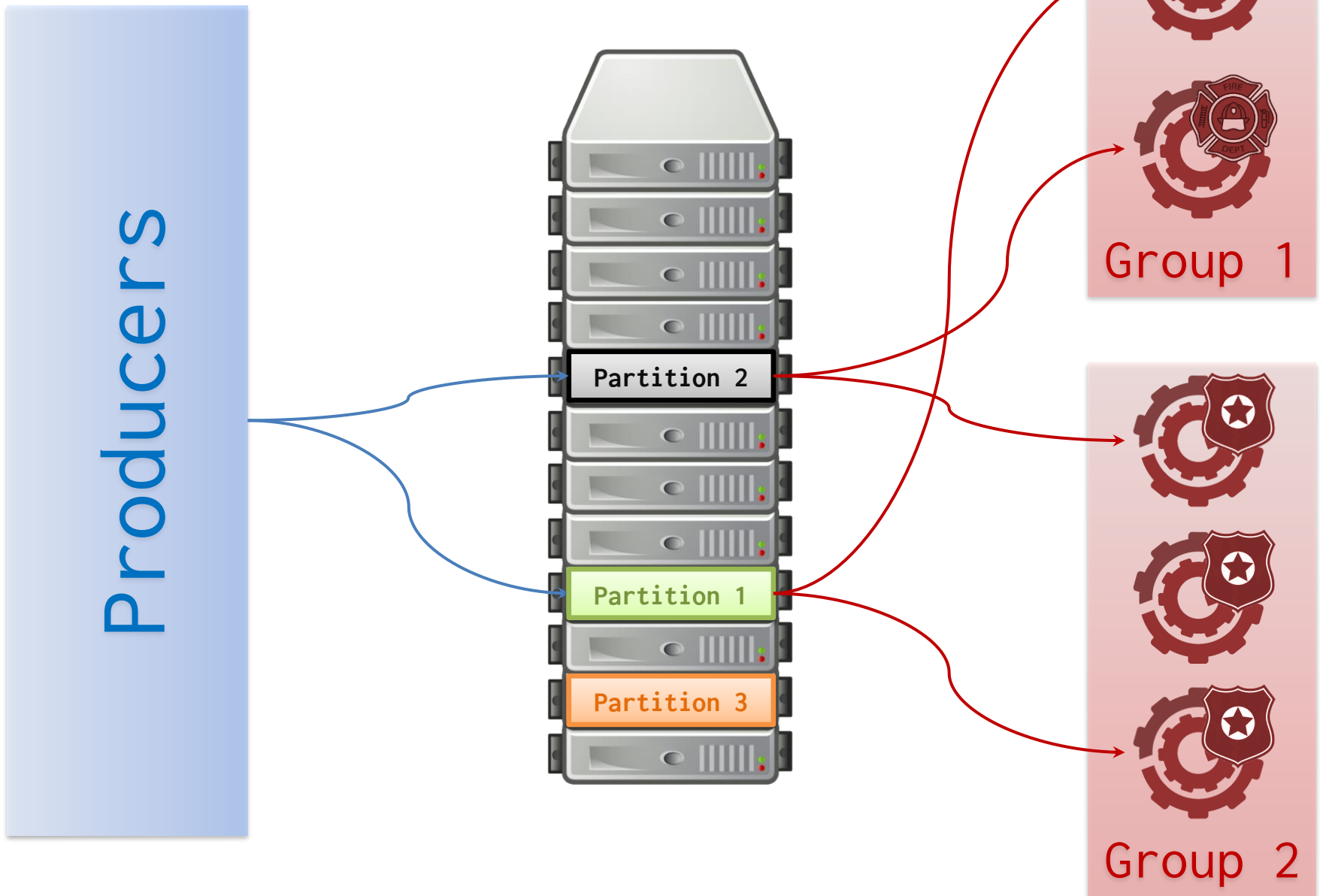
- 1
- 2
- 3
- 4
- 5
- 6

Leader Replication and Reads



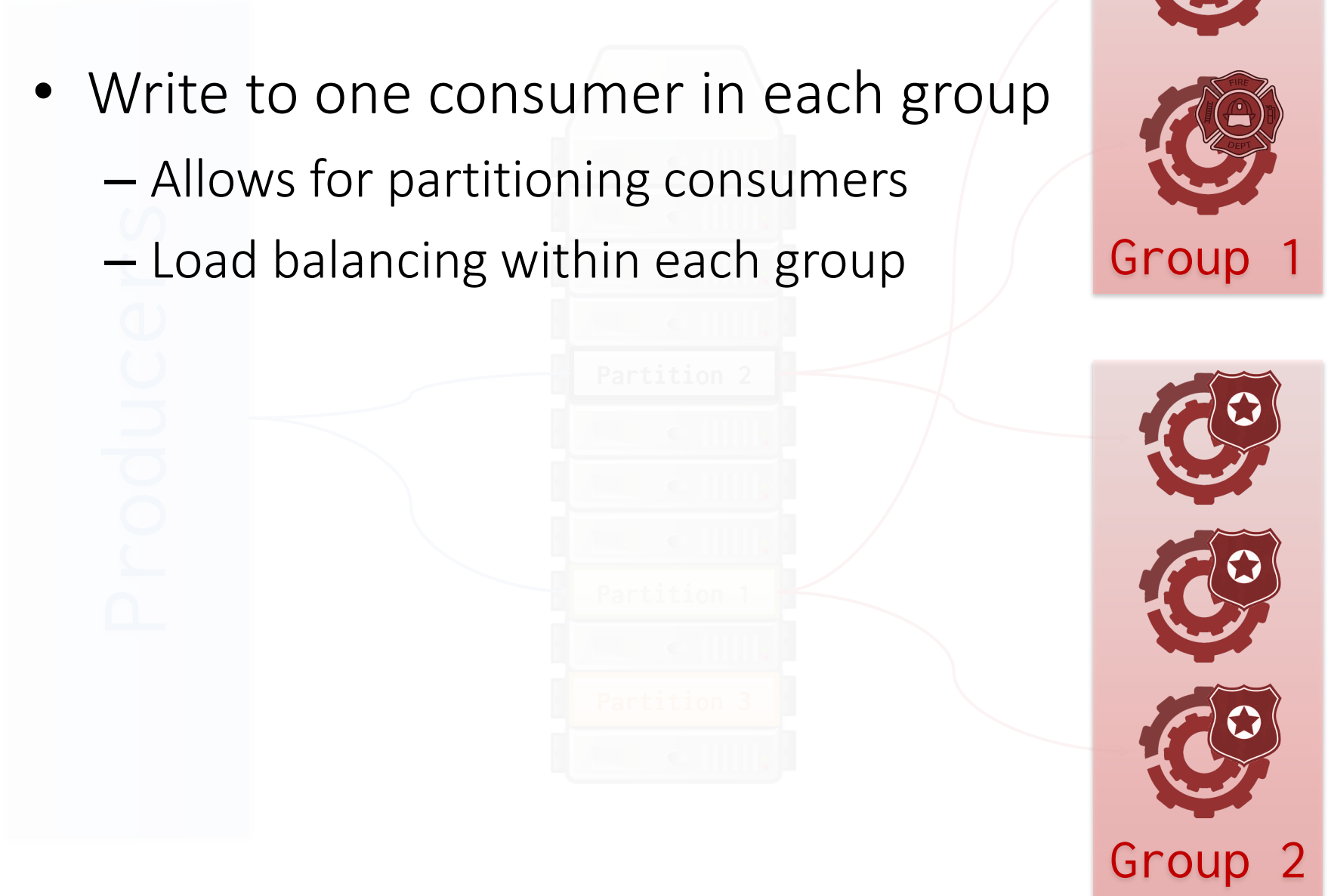
KAFKA: CONSUMER GROUPS

Consumer Groups



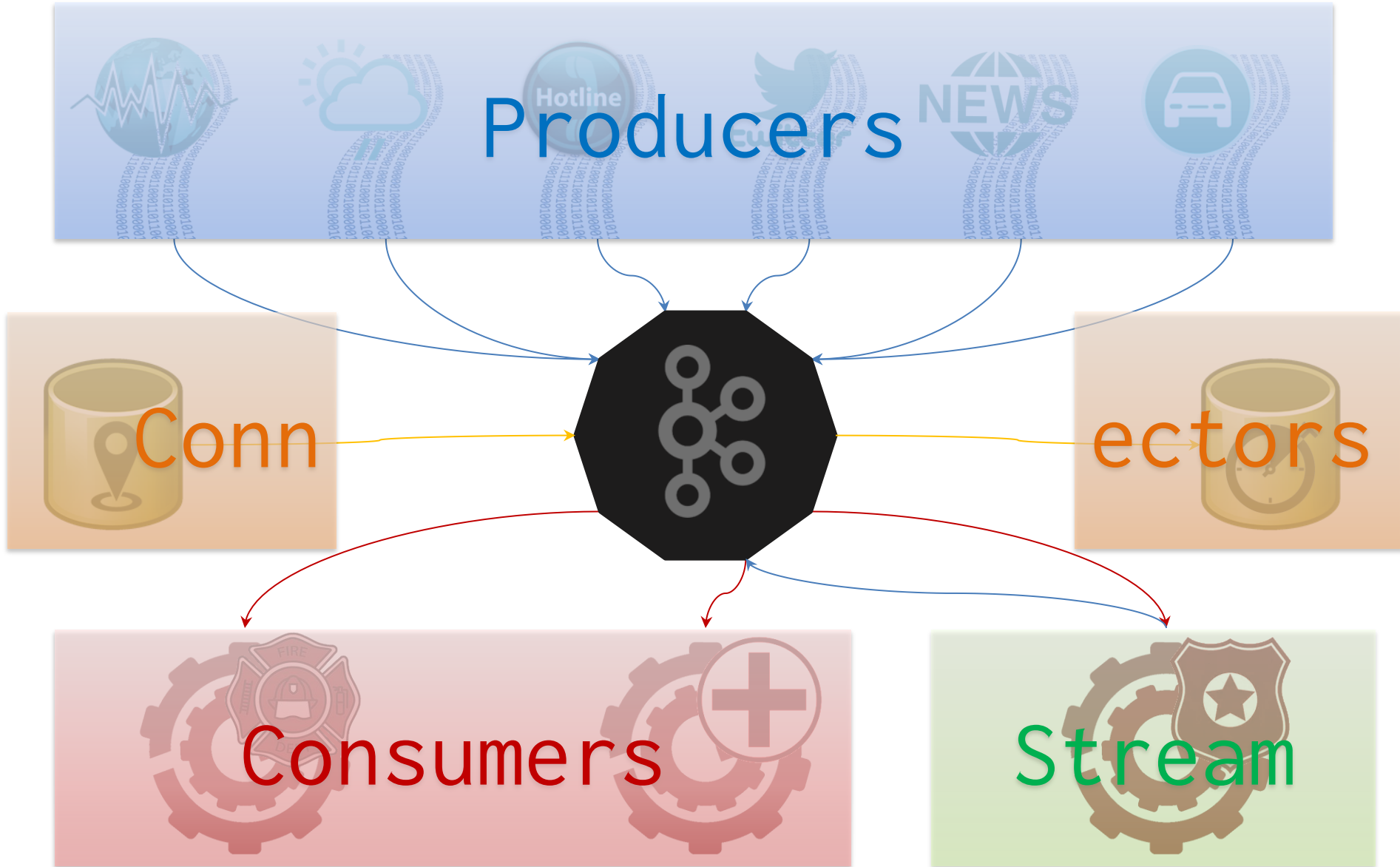
Consumer Groups

- Write to one consumer in each group
 - Allows for partitioning consumers
 - Load balancing within each group



KAFKA: STREAMS AND CONNECTORS

Kafka Overview



Kafka Overview

- **Producer API:**
 - Append records to topics (push)
- **Consumer API:**
 - Read records from topics (pull)
- **Connector API:**
 - Read/write to external components
 - For example, a database or other streaming platforms
- **Stream API (Producer + Consumer):**
 - Read records from input topics
 - Append records to output topics

OPTIMISATIONS AND OTHER FEATURES

Kafka Optimisations

- Log Compaction
 - Repeated sequential values are suppressed
- Direct Disk-to-Network
 - When data don't need to be loaded into JVM
- Consumer / Producer Quotas
 - Set limits to avoid saturating the system
- ...

Kafka Streams API

- **Aggregation** (e.g., count messages)
- **Joins** (e.g., "unify" two streams)
- **Windowing** (define retention period)
- **Continuous Querying** (KSQL)

Available Frameworks





Questions?

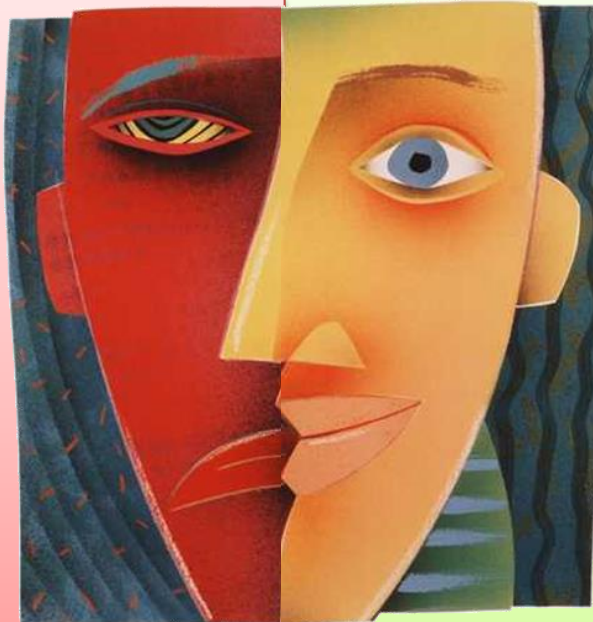
CLASS PROJECTS

Course Marking

- 80% for Weekly Labs
 - 11 labs total
 - Best 9 labs count
- 20% for Class Project

Assignments each week

Working in groups



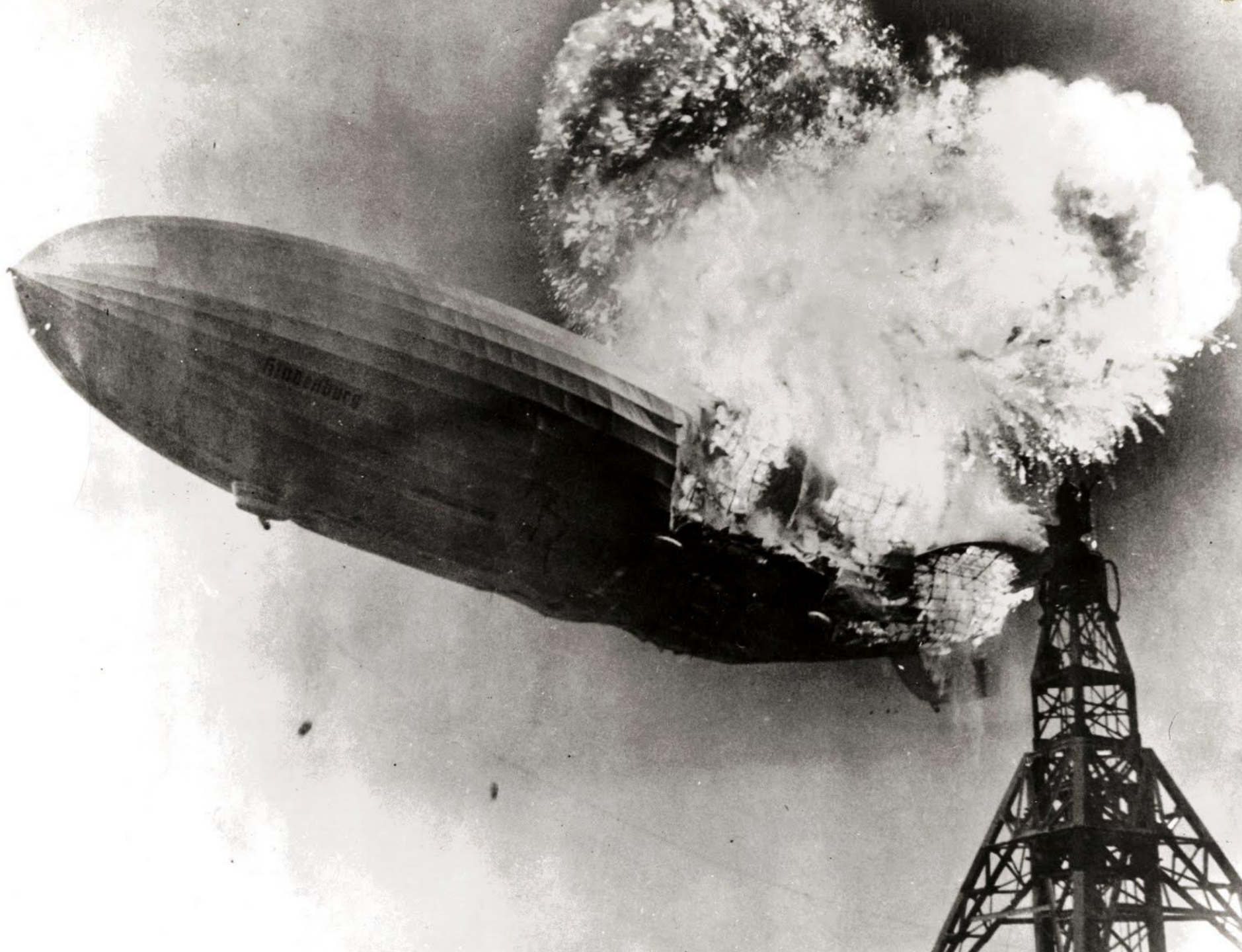
Hands-on each week!

Working in groups!

Class Project



- Done in threes
- Goal: Use what you've learned to do something cool/fun (hopefully)
- Process:
 - We will assign groups tomorrow
 - Start thinking up topics / find interesting datasets!
 - Register topic
 - Work on projects during semester
- Deliverables: 4 minute presentation (video), code repository with documentation in README
- Marked on: Difficulty, appropriateness, scale, good use of techniques, presentation, coolness, creativity, value
 - Ambition is appreciated, even if you don't succeed



Desiderata for project

- Must focus around some technique from the course!
- Expected difficulty: similar to a lab, but without any instructions
- Data not too small:
 - Should have >1,000,000 tuples/entries
- Data not too large:
 - Should have <250,000,000 tuples/entries
 - If very large, perhaps take a sample?

Where to find/explore data?

- Kaggle:
 - <https://www.kaggle.com/>
- Google Dataset Search:
 - <https://datasetsearch.research.google.com/>
- Datos Abiertos de Chile:
 - <https://datos.gob.cl/>
 - <https://es.datachile.io/>
- ...