# CRDTs Illustrated

Arnout Engelen / @raboofje / Xebia
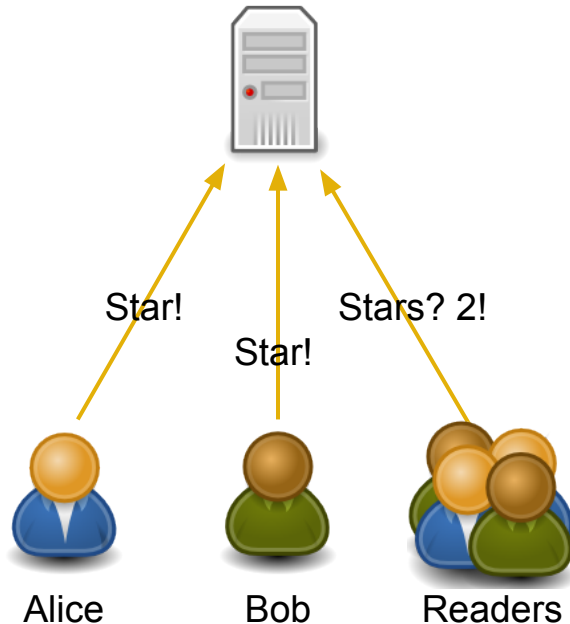
# Today's topic

# Goal

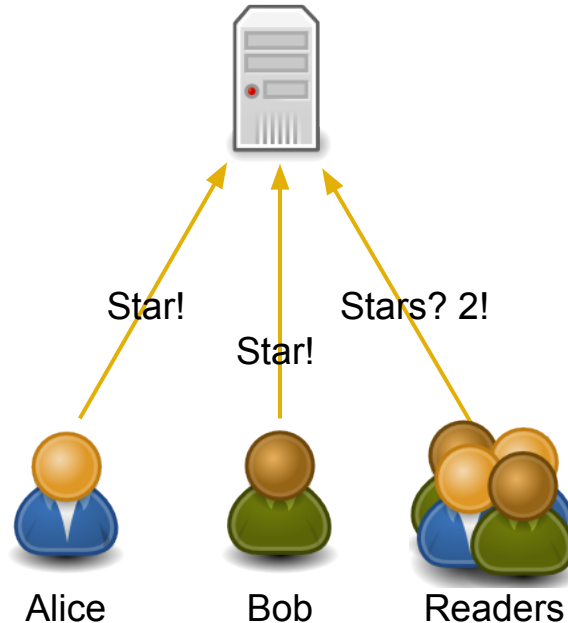Consistency in Distributed Systems

Xebia

# Consistency in Distributed Systems

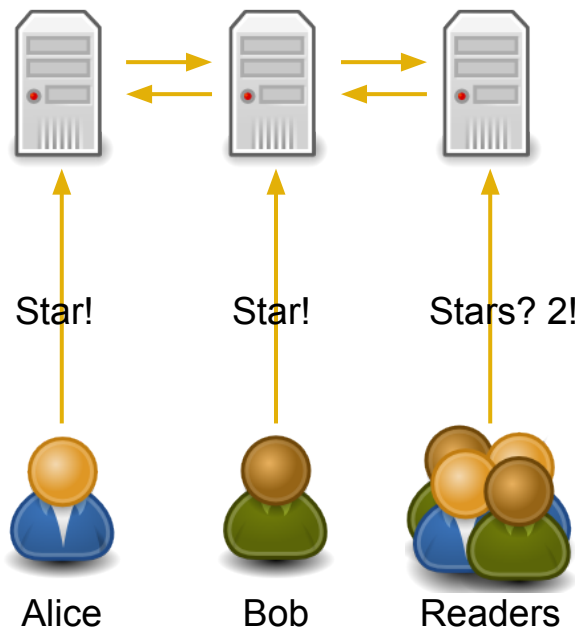Simple (non-distributed) case: single server

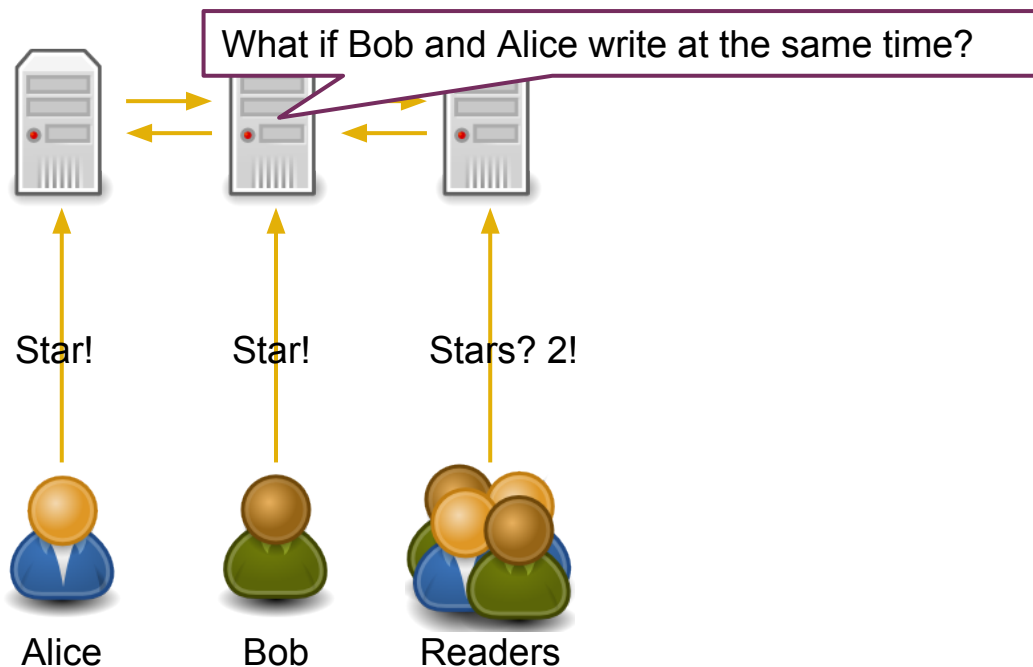# Consistency in Distributed Systems

Problem: network failures :(

# Consistency in Distributed Systems
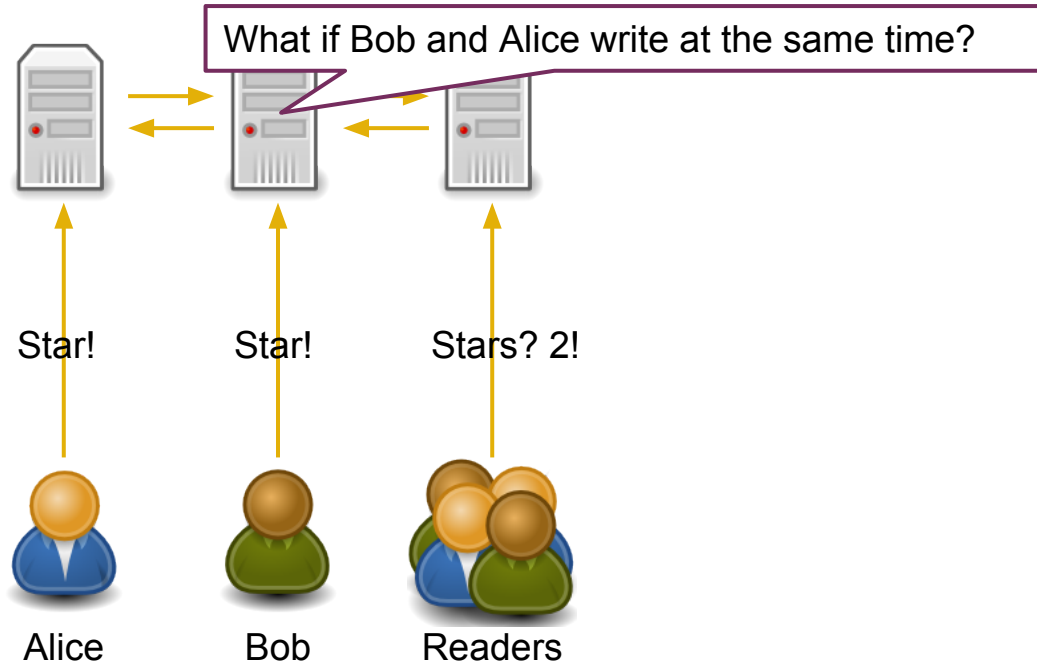
'Solution': distribute!

# Consistency in Distributed Systems

# Consistency in Distributed Systems

Problem: still network failures :(

# Consistency in Distributed Systems

Writing in the presence of failure

- Strong Consistency
  - Sequential writes
  - Impossible when A and B are disconnected
  - "No availability in case of network partitions"

Xebia

# Consistency in Distributed Systems

Writing in the presence of failure

- Strong Consistency
  - Sequential writes
  - Impossible when A and B are disconnected
  - "No availability in case of network partitions"

- Eventual Consistency
  - Update partitions independently, converge 'eventually'
  - Complicated algorithms, hard to verify/test

Xebia

# CRDTs

# CRDTs

Strong Eventual Consistency

- Once you've seen the same events, you're (immediately) in the same state

Xebia

# CRDT: counter

action: *plus 5*
value: 0       5

action:
value: 0           5

Xebia

# CRDT: counter

action: *plus 5*   *minus 4*

value: 0   5   1   -2

action:   *minus 3*

value: 0   5   2   -2

# CRDT: Op-based counter

# CRDT: Op-based counter



action: *plus 5*    *minus 4*         *times 2*
value: 0      5         1      -2       -4    -5

action:        *minus 3*       *minus 1*
value: 0         5       2    -2    -3    -6

# CRDT: Op-based counter

| action: | *plus 5* | *minus 4* | | *times 2* | | |
|---|---|---|---|---|---|---|
| value: 0 | 5 | 1 | -2 | -4 | -5 | |

| action: | | *minus 3* | | *minus 1* | | |
|---|---|---|---|---|---|---|
| value: 0 | | 5 | 2 | -2 | -3 | -6 |

Op-based CRDT rule 1:

- All concurrent operations must commute
  - (x-4)-3 == (x-3)-4
  - (x*2)-1 != (x-1)*2

Xebia

# CRDT: Op-based counter

action: *plus 5*
value: 0      5

action:
value: 0      5   10

# CRDT: Op-based counter

action: *plus 5*
value: 0        5

action:
value: 0              5    10

Op-based CRDT rule 2:

- Updates must be applied exactly once
  - Applying '+5' twice invalidates the result

Xebia

# Exactly Once Delivery

# Exactly Once Delivery

- .. is generally impossible when partitions happen
  - Who acknowledges the acknowledgement?

Xebia

# Exactly Once Delivery

- .. is generally impossible when partitions happen
  - Who acknowledges the acknowledgement?

- Pick one:
  - At Most Once Delivery (fire & forget)
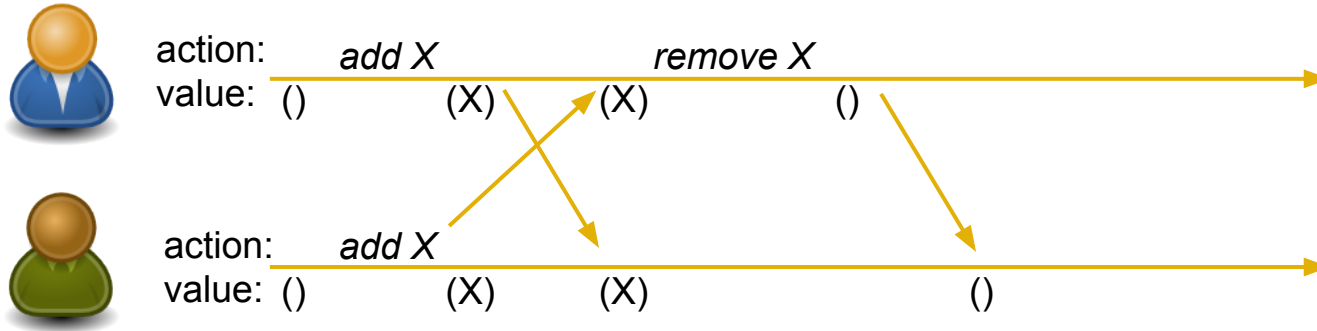  - At Least Once Delivery (retry)

Xebia

# Exactly Once Delivery

- .. is generally impossible when partitions happen
  - Who acknowledges the acknowledgement?

- But: Exactly Once Delivery *Semantics* are possible!
  - when processing the same message again has no effect
  - *Idempotence*

Xebia

# Sets: naive approach

action: *add X*                           *remove X*

value: ()          (X)        (X)              ()

action: *add X*

value: ()         (X)        (X)                   ()

# Sets: naive approach

action: *add X*       *remove X*

value: ()    (X)        ()     ???

action: *add X*

value: ()    (X)    (X)        ???

Xebia

# CRDT: Observed-Remove Set

# CRDT: Observed-Remove Set

# CRDT: Observed-Remove Set



a — action: *add $X_a$* ... *remove $X_a$*
value: () ($X_a$) () ()

b — action:
value: () () ($X_a$)

Op-based CRDT rule 3:

- Updates must be applied in-order
  - (in which they were sent from their origin)

# CRDTs

Operation-based CRDTs:

- All (concurrent) operations must *commute*
- Require *exactly-once* delivery semantics
- Require *in-order* delivery

Xebia

# CRDTs

Operation-based CRDTs:

- All (concurrent) operations must *commute*
- Require *exactly-once* delivery semantics
- Require *in-order* delivery

Next up: State-based CRDTs

Xebia

# CRDT: State-based counter



a

action:    *plus 5*
value:  {}=0      {a:5}=5

b

action:
value:  {}=0         {a:5}=5

# CRDT: State-based counter

a

action: *plus 5*
value: {}=0     {a:5}=5

b

action:
value: {}=0     {a:5}=5

State-based CRDT:

- Local *update*
- Send state and *merge*

Xebia

# CRDT: State-based counter

a

action:     *plus 5*          *plus 4*

value:   {}=0     {a:5}=5       {a:9}=9     {a:9,b:3}=12

b

action:

value:   {}=0     {a:5}=5    *plus 3*   {a:5,b:3}=8    {a:9,b:3}=12

State-based CRDT:

- Local *update*
- Send state and *merge*

# CRDT: State-based counter

a

action:
value: {}=0    *plus 5* {a:5}=5    *plus 4*    {a:9}=9    {a:9,b:3}=12

b

action:
value: {}=0    {a:5}=5    *plus 3*    {a:5,b:3}=8    {a:9,b:3}=12

State-based CRDT rule 1:

- We allow retransmissions.
  - The merge function should be *idempotent*.

# CRDT: State-based counter

a

action:      *plus 5*              *plus 4*
value:  {}=0          {a:5}=5          {a:9}=9

In-order delivery not assumed!

b

action:
value:  {}=0                    {a:9}=9      ???

Xebia

# CRDT: State-based counter



**State-based CRDT rule 2:**

- Output must be independent of the order of merges
  - The merge function is *commutative* and *associative*

# CRDT: State-based counter

a

action:        *plus 5*         *plus 4*
value:   {}=0      {a:5}=5       {a:9}=9

In-order delivery not assumed!

b

action:
value:   {}=0                           {a:9}=9     {a:9}=9

Our merge function here takes the 'max':

This is an *increment-only counter*.

Xebia

# CRDT: State-based counter

a

action:       *plus 5*       *minus 4*       *plus 3*

value:   {}=0      {a:+5}=5      {a:+5,-4}=1      {a:+8,-4}=4

b

action:

value:   {}=0              {a:+5,-4}=1   {a:+5,-4}=1

Solution 1: combine 2 counters

(positive and negative: 'PN-counter')

Xebia

# CRDT: State-based counter

a

action:     *plus 5*            *minus 4*            *plus 3*
value:  {}=0        {a:5}=5         {a':1}=1          {a'':4}=4

b

action:
value:  {}=0                          {a':1}=1     {a':1}=1

Solution 2: version vectors

Xebia

# CRDT: State-based counter

a

action:   *plus 5*                    *plus 4*

value:  {}=0       {a:5}=5                 {a:9}=9       {a:9,b:3}=12

b

action:                        *plus 3*

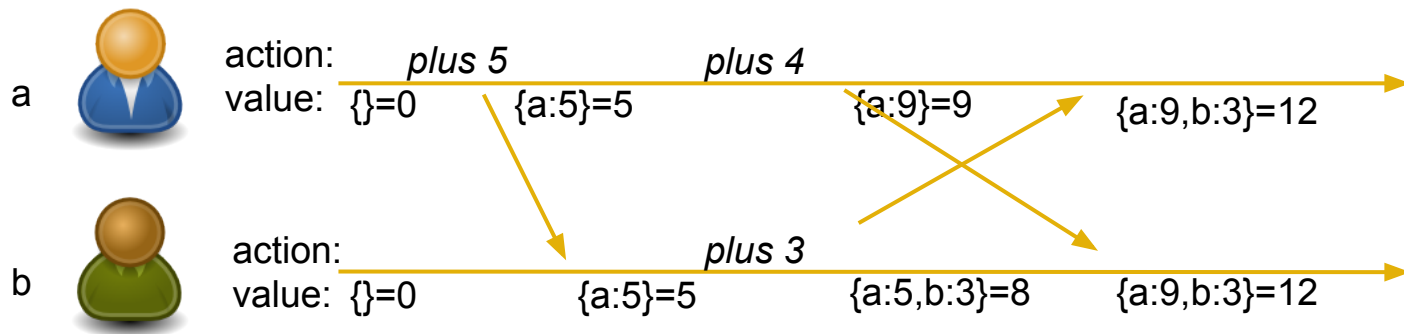value:  {}=0          {a:5}=5       {a:5,b:3}=8   {a:9,b:3}=12

State-based CRDT rules 3 and 4:

- We need a concept of 'going forward' (growing, clocks)
- Updates and merges always go forward
  - Concurrent states may not be comparable ({a:9} and {a:5,b:3})

Xebia

# CRDT: State-based counter

a

action: *plus 5*       *plus 4*

value: {}=0    {a:5}=5       {a:9}=9      {a:9,b:3}=12

b

action:             *plus 3*

value: {}=0      {a:5}=5    {a:5,b:3}=8    {a:9,b:3}=12

State-based CRDT rules 3 and 4:

- We need a concept of 'going forward' (growing, clocks)
- Updates and merges always go forward

In other words:

- There is a *partial order* on states
- Updates and merges must *increase* the state in this order

**Xebia**

# CRDTs: 2 kinds, 2 sets of rules

Operation-based CRDTs (Commutative, CmRDTs):

- All (concurrent) operations must be *commutative*
- Require *unique* and *in-order* delivery

State-based CRDTs (Convergent, CvRDTs):

- merge must be *idempotent*
- merge must be *commutative* and *associative*
- there exists a *partial order* on the states
- merge and update both *increase* the state along this order

# DEMO

Xebia

# Conclusions

- CRDTs:
  - Have *simple rules*
  - Guarantee *Strong Eventual Consistency*
  - Can be *composed* to build more complex structures

  - Cannot model *every* data type
  - .. see if you can 'bend your problem'

Xebia

# Further Material

- Watch Marc Shapiro's talks
  - (e.g. Strong Eventual Consistency at MS Research)

- Check out @cmeik's reading list
  - http://christophermeiklejohn.com/crdt/2014/07/22/readings-in-crdts.html

- Check out @cmeik's talk after lunch (Theater)

Xebia

# Questions?

# Thank you!

-- Arnout Engelen, @raboofje

Xebia