# A principled approach to programming distributed systems

Marc Shapiro

*Summary School on Distributed and Replicated Environments*
*DARE 2023*

SORBONNE UNIVERSITÉ CRÉATEURS DE FUTURS DEPUIS 1257  LIP

*informatiques* *mathématiques*
Inria

---

## CRDTs: not a silver bullet

CRDT:
- Extend sequential data type (commutative or not)
- Merge concurrent updates: commutative, associative, idempotent
- Guarantee availability + convergence

Designing distributed applications remains challenging
- Multiple CRDTs
- Correct interactions?

Experiments: calendar, student registration, game tournaments, auction, ticket reservations, file system, etc.

Real applications: LWW, Bet365, text editor, ???

---

## CAP trade-offs

Strong consistency
- High system cost
  - One-off
- Low programmer cost
- Sequential bottleneck
- partition $\implies$ ¬available

Too strong: slow

"Weak" consistency
- Responsive
- High programmer cost
  - Recurring
- Asynchronous
- Available

Too weak: buggy

"As fast as possible, as strong as necessary"

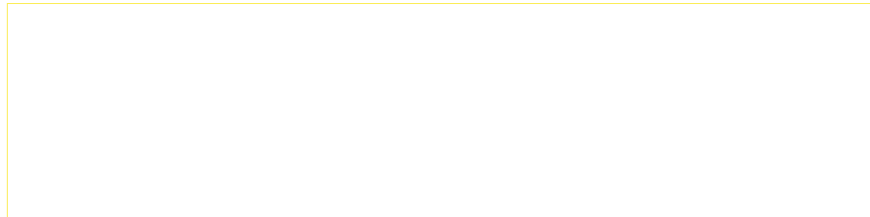Strong enough; no stronger

---

## Consistent ≜ maintains invariant

Invariant ≜ predicate over the state of the system
- Credit cards: *hash(n) = 0*
- Seat reservations: *remaining_seats ≥ 0*
- Social network: *friend (A,B) ⟺ friend (B,A)*
- Bank: *Σ accounts = constant*
- File system: *tree ∧ …*
- Student registration:
  *attends (student, course) ⟹ registered (student)*
- Storage backend: *journal.low ≤ checkpoint.high*

Safe system ≜ invariant is true in every observable state
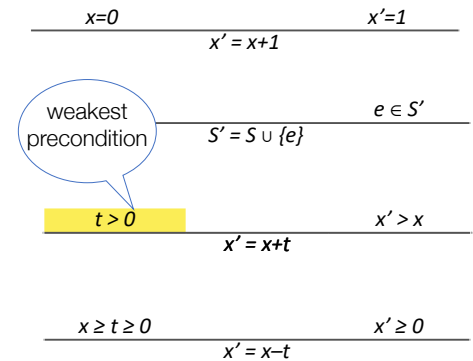
System guarantee vs. application-level effort

# Sequential safety

# Precondition, postcondition

$$\frac{Pre(\sigma) \qquad\qquad Post(\sigma')}{\sigma \;—\; U \;\longrightarrow\; \sigma'}$$
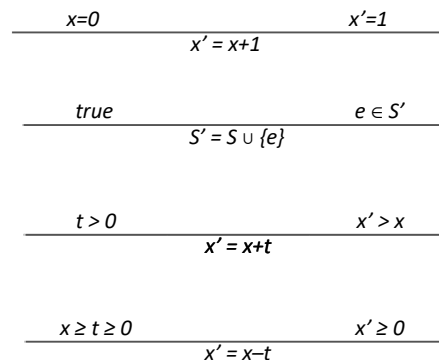
Update $U$

Transitions the state from $\sigma$ to $\sigma'$

Assuming $Pre(\sigma)=true$

    then $U$ ensures $Post(\sigma')$

$$\frac{x=0 \qquad\qquad x'=1}{x' = x+1}$$

weakest precondition

$$\frac{e \in S'}{S' = S \cup \{e\}}$$

$$\frac{t > 0 \qquad\qquad x' > x}{x' = x+t}$$

$$\frac{x \geq t \geq 0 \qquad\qquad x' \geq 0}{x' = x-t}$$

# Precondition, postcondition

$$\frac{Pre(\sigma) \qquad\qquad Post(\sigma')}{\sigma \;—\; U \;\longrightarrow\; \sigma'}$$

Update $U$

Transitions the state from $\sigma$ to $\sigma'$

To ensure $Post(\sigma')$,

    it must be that $Pre(\sigma)=true$

$$\frac{x=0 \qquad\qquad x'=1}{x' = x+1}$$

$$\frac{true \qquad\qquad e \in S'}{S' = S \cup \{e\}}$$

$$\frac{t > 0 \qquad\qquad x' > x}{x' = x+t}$$

$$\frac{x \geq t \geq 0 \qquad\qquad x' \geq 0}{x' = x-t}$$

# Ex.: Generate unique identifier

UIDs      [ ]      [0,0]    [0,1]    [0,2]    [0,3]     *time*

new_uid()  new_uid()  new_uid()  new_uid()   $\longrightarrow$
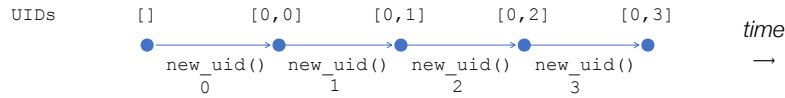
0         1         2         3

$$\frac{M \text{ is a set} \qquad n \notin M \qquad M' \text{ is a set}}{M' = M \cup \{n\}}$$

```
Inv() = { IDs = [0, last] }
static last = -1;
new_uid() = {
  precondition Inv();     // multiset is a set, max=last
  last++;                 // ∉ IDs
  return last;            // is unique
  postcondition Inv();
}
```

# ∀ update safe ⟹ seq. system safe

UIDs  []        [0,0]      [0,1]      [0,2]      [0,3]

●————→●————→●————→●————→●        *time*
new_uid()  new_uid()  new_uid()  new_uid()        →
0          1          2          3

Individually
Safe
"C"

Safe state ≜ satisfies invariant

In a sequential system, if:
- The initial state satisfies *Inv*
- Every update *U* has
  precondition *wp(U, Inv)*

Then the system is safe

$$\forall U \quad \frac{Inv(\sigma) \quad wp(U, Inv) \quad Inv(\sigma')}{\sigma \; —\!U\!—\!\!\longrightarrow \sigma'} \quad \Longrightarrow \text{safe}$$

---

# Examples

Issue credit card

$$\frac{hash\,(n) = 0}{S' = S \cup \{n\}}$$

Deliver medication
- *Inv(r) { r ≥ 0}*        *// r = remaining number of boxes*

Increase
allocation:
$$\frac{Inv(r) \qquad\qquad\qquad Inv(r')}{r' = r+1}$$

Deliver a box:
$$\frac{Inv(r) \qquad r \geq 1 \qquad Inv(r')}{r' = r-1}$$

---

# Ex.: maintain checkpoint

checkpoint        *hc*

·········|————————————|

          |————————————·········
          *lj*        journal

*Inv(σ) = { σ.lj ≤ σ.hc}*        *// no gap*

Advance
checkpoint:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad\qquad Inv(\sigma')}{\sigma'.hc = \sigma.hc + d}$$

Advance
journal:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad \sigma.lj + d \leq \sigma.hc \qquad Inv(\sigma')}{\sigma'.lj = \sigma.lj + d}$$

Advance both:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad Inv(\sigma'')}{\sigma'.hc = \sigma.hc + d; \; \sigma''.lj = \sigma'.lj + d}$$

---

# Ex.: maintain checkpoint

checkpoint        *hc*

·········|————————————|

               |————————————·········
               *lj*        journal

*Inv(σ) = { σ.lj ≤ σ.hc}*        *// no gap*

Advance
checkpoint:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad\qquad Inv(\sigma')}{\sigma'.hc = \sigma.hc + d}$$

Advance
journal:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad \sigma.lj + d \leq \sigma.hc \qquad Inv(\sigma')}{\sigma'.lj = \sigma.lj + d}$$

Advance both:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad Inv(\sigma'')}{\sigma'.hc = \sigma.hc + d; \; \sigma''.lj = \sigma'.lj + d}$$

# Ex.: maintain checkpoint



$Inv(\sigma) = \{ \sigma.lj \leq \sigma.hc\}$     // no gap

Advance checkpoint:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad Inv(\sigma')}{\sigma'.hc = \sigma.hc + d}$$

Advance journal:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad \sigma.lj + d \leq \sigma.hc \qquad Inv(\sigma')}{\sigma'.lj = \sigma.lj + d}$$

Advance both:
$$\frac{Inv(\sigma) \wedge d \geq 0 \qquad\qquad Inv(\sigma'')}{\sigma'.hc = \sigma.hc + d; \; \sigma''.lj = \sigma'.lj + d}$$

# File system: tree invariant + mv



mv(B,C)             **mv(A,B)**

$$\frac{Inv(\sigma) \qquad \neg \sigma.path(c, np) \qquad Inv(\sigma')}{\sigma — mv\,(c, np) \longrightarrow \sigma'}$$

$Inv(\sigma) =$    // Tree invariant
- Every node is reachable from root
- A node has a single parent (but root has none)
- Names unique per directory
- Acyclic graph

# Sequential safety: summary

Safe state ≜ satisfies invariant *Inv*

In a sequential system, if:
- The initial state is safe
- Every update *U* is individually safe
  - i.e., has a precondition *Pre(U)* such that *Inv* remains true after *U*
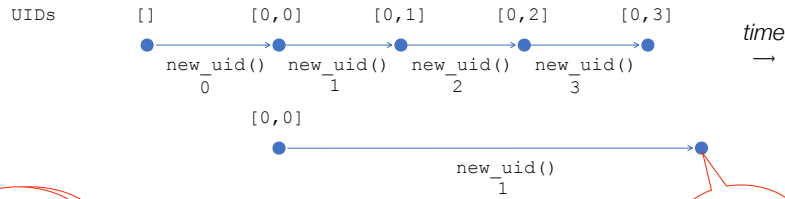  - i.e., $Pre(U) \implies wp(U, Inv)$

then every state of the system is safe

# Concurrency anomalies

# Concurrent generate UID

UIDs
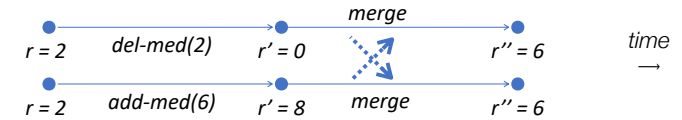
```
Inv() = { IDs = [0, last] }
static last = -1;
new_uid() = {
  precondition Inv();    // multiset is a set, max=last
  last++;                // ∉ IDs
  return last;           // is unique
  postcondition Inv();
}
```
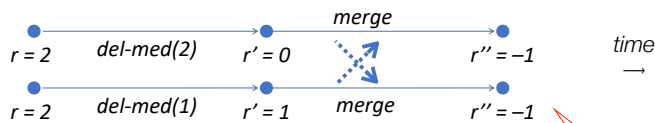
TOCTTOU

Anomaly!

[0,0] [0,1] [0,2] [0,3]

new_uid() 0 new_uid() 1 new_uid() 2 new_uid() 3

[0,0]

new_uid() 1

*time* →

# Concurrent delivery of medications (1)

$r = 2$   *del-med(2)*   $r' = 0$   *merge*   $r'' = 6$

$r = 2$   *add-med(6)*   $r' = 8$   *merge*   $r'' = 6$

*time* →

$$\text{add-med}(n) \quad \frac{r \geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$\text{del-med}(n) \quad \frac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$
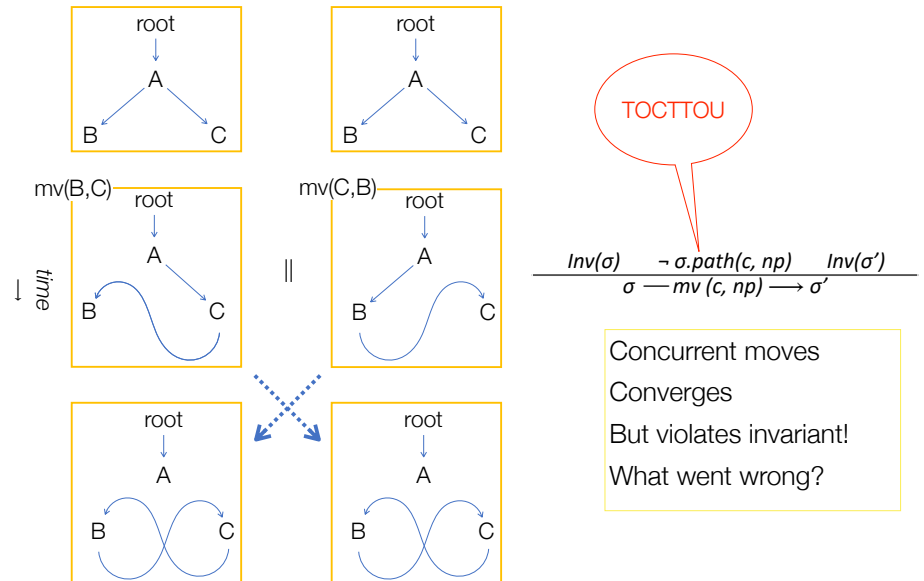
harmless TOCTTOU

# Concurrent delivery of medications (2)

$r = 2$   *del-med(2)*   $r' = 0$   *merge*   $r'' = -1$

$r = 2$   *del-med(1)*   $r' = 1$   *merge*   $r'' = -1$

*time* →

harmful TOCTTOU

Anomaly!

$$\text{add-} \quad \frac{\geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$\text{del-med}(n) \quad \frac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$

# FS: Concurrent move anomaly

root → A; B ← A → C

root → A; B ← A → C

mv(B,C)   root → A; B, C

||

mv(C,B)   root → A; B, C

↓ *time*

TOCTTOU

root → A; B, C

root → A; B, C

$$\frac{Inv(\sigma) \qquad \neg\, \sigma.path(c, np) \qquad Inv(\sigma')}{\sigma \longrightarrow mv\,(c, np) \longrightarrow \sigma'}$$

Concurrent moves
Converges
But violates invariant!
What went wrong?

# File system: the DFS-R bug

2003: DFS-R: Windows NTFS replication layer
- large industrial customers
- unexplained data loss

2007: model checking exposes move anomaly: completely unexpected

2021: continues to bite developers
- Google Drive diverges
- Dropbox duplicates

Reasoning about concurrency is hard!

# Concurrent anomalies: summary

Concurrency allows anomalies
- Incorrect behavior ≜ violates invariant
- Violation does not occur in a sequential execution

Anomalies seem to be caused by TOCTTOU race on precondition

Design objectives:
- Minimise remote coordination
- Allow benign TOCTTOUs
- Avoid harmful TOCTTOUs

But concurrency reasoning is too hard!

# CISE: a sound approach to safe distributed systems

# 'Cause I'm Strong Enough (CISE)

Distributed system with CRDTs

Systematic method to prove whether two updates may execute concurrently without harm
- Concurrency not provably harmless is assumed harmful

Fully formalized, proven sound
- Generalises sequential correctness
- An application of the rely-guarantee logic

My presentation is informal

Gotsman et al., *'Cause I'm strong enough: Reasoning about consistency choices in distributed systems*, POPL 2016, DOI 10.1145/2837614.2837625

## CISE conditions

For all updates *U, U':*

(1) Sequentially safe:
- Initial state satisfies invariant *Inv*
- Precondition of *U* satisfies the weakest-precondition of the invariant *wp(U,Inv)*
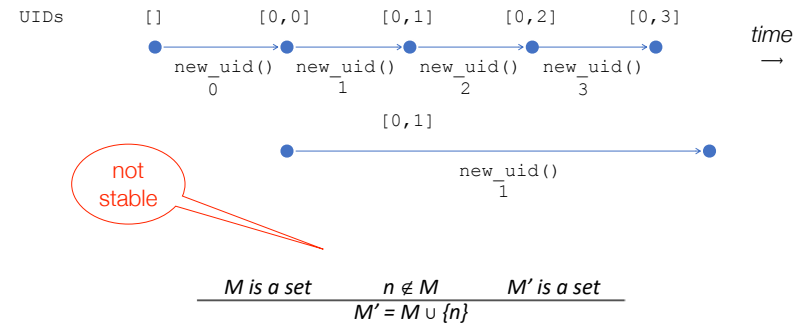
> Individually Safe

(2) Convergent:
- When *U* concurrent to *U'*
- replace *U* ∥ *U'* with *merge(U,U')*
- *merge(U,U')* commutative, associative, idempotent
- *merge(U,U')* preserves *Inv*

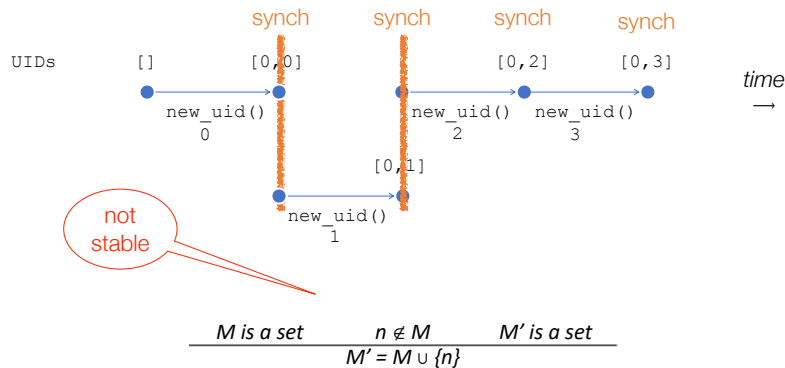> harmless TOCTTOU

(3) Stable precondition:
- When *U* concurrent to *U'*
- the precondition of *U* is not made false by *U'*
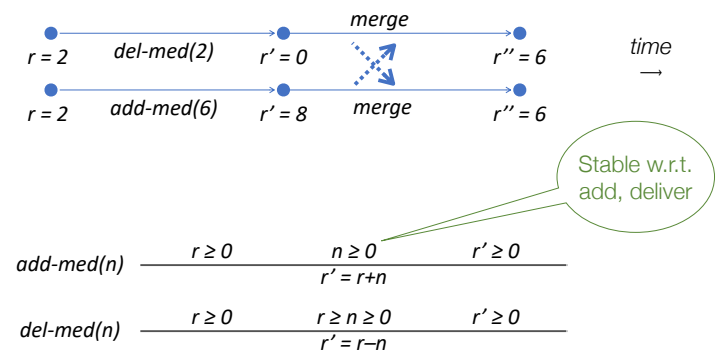
## CISE Concurrent generate UID



> not stable

$$\frac{M \text{ is a set} \quad n \notin M \quad M' \text{ is a set}}{M' = M \cup \{n\}}$$

Solution: serialise ⟹ consensus

## CISE Concurrent generate UID



> not stable

$$\frac{M \text{ is a set} \quad n \notin M \quad M' \text{ is a set}}{M' = M \cup \{n\}}$$

Solution: serialise ⟹ consensus

## CISE Concurrent delivery of medications (1)



> Stable w.r.t. add, deliver

$$add\text{-}med(n) \quad \frac{r \geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$del\text{-}med(n) \quad \frac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$

- *add-med* ∥ *add-med* OK
- *del-med* ∥ *add-med* OK

# CISE Concurrent delivery of medications (2)



*time* →

$$add\text{-}med(n) \quad \dfrac{r \geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$del\text{-}med(n) \quad \dfrac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$

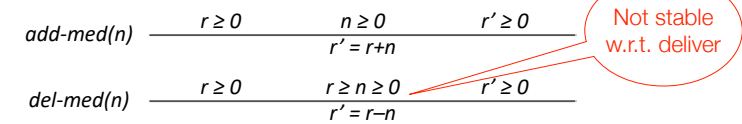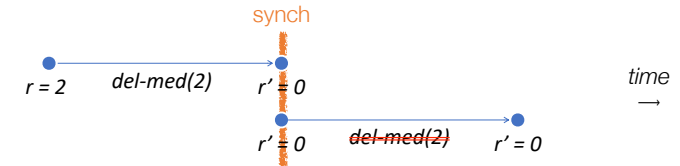**Not stable w.r.t. deliver**

*del-med* design alternatives:
- Allow over-delivery → possibly punish after the fact
- Synchronise *del-med; del-med* → lose availability
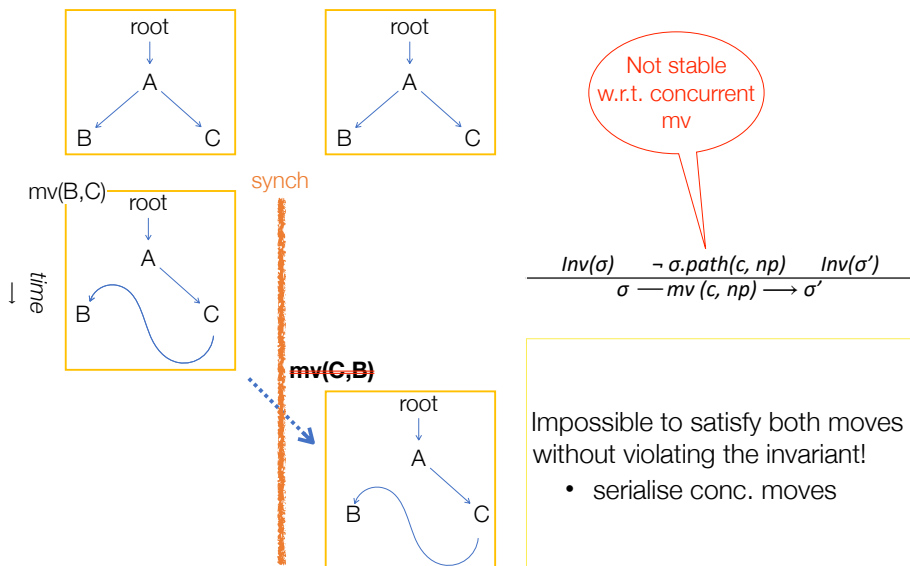
**weaken invariant, compensate**
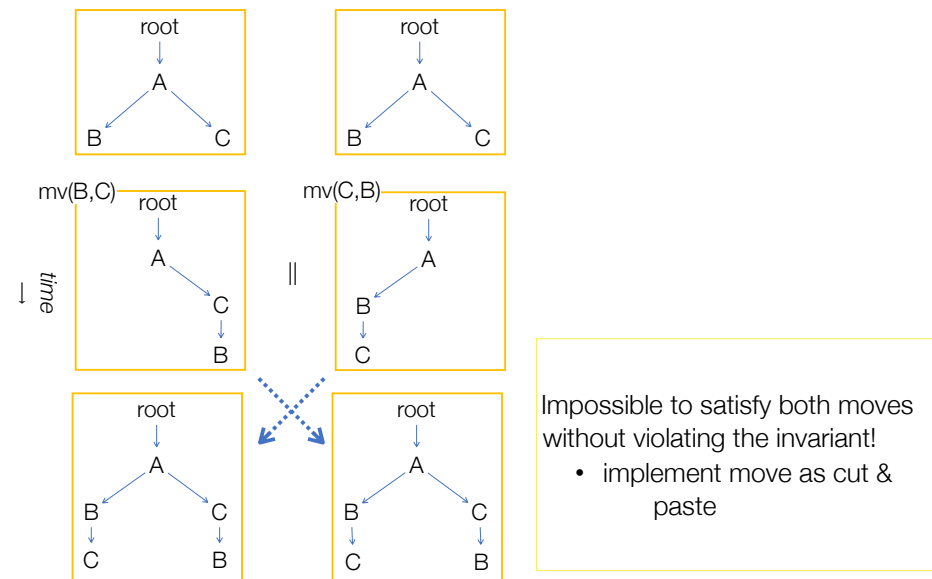
# CISE Concurrent delivery of medications (3)

synch



*time* →

$$add\text{-}med(n) \quad \dfrac{r \geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$del\text{-}med(n) \quad \dfrac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$

**Not stable w.r.t. deliver**

*del-med* design alternatives:
- Allow over-delivery → possibly punish after the fact
- Synchronise *del-med; del-med* → lose availability

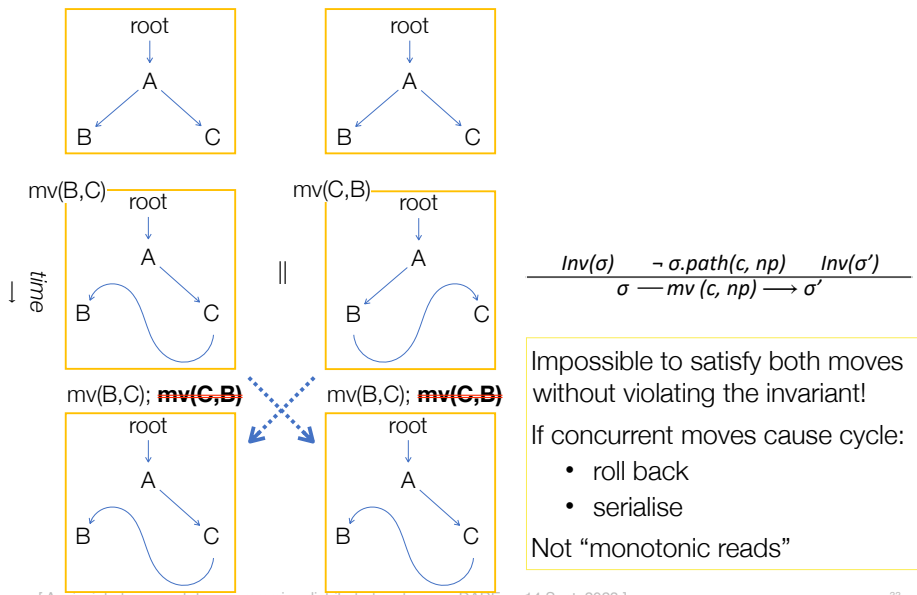# CISE Tree: serialise move



**Not stable w.r.t. concurrent mv**

$$\dfrac{Inv(\sigma) \qquad \neg\, \sigma.path(c, np) \qquad Inv(\sigma')}{\sigma \longrightarrow mv\,(c, np) \longrightarrow \sigma'}$$

Impossible to satisfy both moves without violating the invariant!
- serialise conc. moves

# Dropbox: *mv* as cut & paste



Impossible to satisfy both moves without violating the invariant!
- implement move as cut & paste

## Tree + mv: Kleppmann's approach



$$\frac{Inv(\sigma) \quad \neg\, \sigma.path(c, np) \quad Inv(\sigma')}{\sigma \longrightarrow mv\,(c, np) \longrightarrow \sigma'}$$

Impossible to satisfy both moves without violating the invariant!

If concurrent moves cause cycle:
- roll back
- serialise

Not "monotonic reads"

## Summary: Tree + mv

Sequentially correct: sequential moves are OK
- Weakest precondition: not mv under self

CISE: precondition not stable under concurrent *mv*

Known design options:
- no *mv* op (XML)
- no *mv* op, copy-paste ⟹ duplicates (Dropbox)
- *up-mv* vs. *down-mv* [Nair 2021]
  - *up-mv* ‖ *up-mv* stable
  - *up-mv* ‖ *down-mv* stable
  - *down-mv* ‖ *down-mv* not stable
- serialise *a priori:* lock [Najafzadeh 2018]
- serialise *a posteriori*: non-monotonic [Kleppmann 2022]

## CISE conditions (again!)

For all updates *U, U':*

(1) Sequentially safe:
- Initial state satisfies invariant *Inv*
- Precondition of *U* satisfies the weakest-precondition of the invariant *wp(U,Inv)*

(2) Convergent:
- When *U* concurrent to *U'*
- *U ‖ U'* convergent
  - *merge(U,U')* preserves *Inv*

(3) Stable precondition:
- When *U* concurrent to *U'*
- the precondition of *U* is not made false by *U'*

## Using CISE in practice

Manual:
1. Individually correct: manual, testing
2. Convergence: CRDT library
3. Stability: consider all pairs of possibly concurrent updates

Tools
- Library: Bounded Counter [Balegas SRDS 2015]
- Stand-alone: specification language + SMT solver
  - CISE Tool [Najafzadeh 2015]
  - Soteria [Nair ESOP 2020]
  - BLOOM [Alvaro CIDR 2011]
- Integrated language/compiler
  - Conflict-Aware Replicated Data Types [arXiv 1802.08733]
  - LoRe [Haas ECOOP 2023]

## Summary: As fast as possible

*U, U'* individually safe ∧ convergent ∧ mutually stable
- May execute concurrently
- Availability
- Perfect scalability

## Summary: As strong as necessary

*U, U'* must not execute concurrently if:
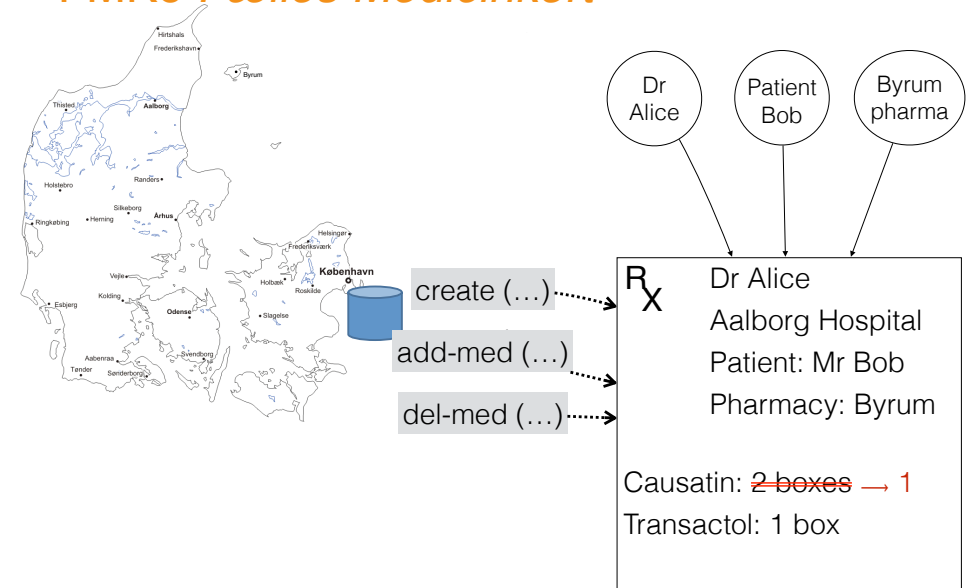
  not commutative ∧ not convergent

 ∨ not mutually stable

Design options:
- Refine invariant (e.g., bank account number)
- Downgrade invariant
  - from *"require x≥0"* to *"prefer x≥0"*
- Compensate: Weaken invariant, repair
- Serialise: *U; U'* or *U'; U*
  - Lock, single server, social convention, etc.
  - Monotonic, a priori: consensus
  - A posteriori: rollback.  Finality?

Analyse again!

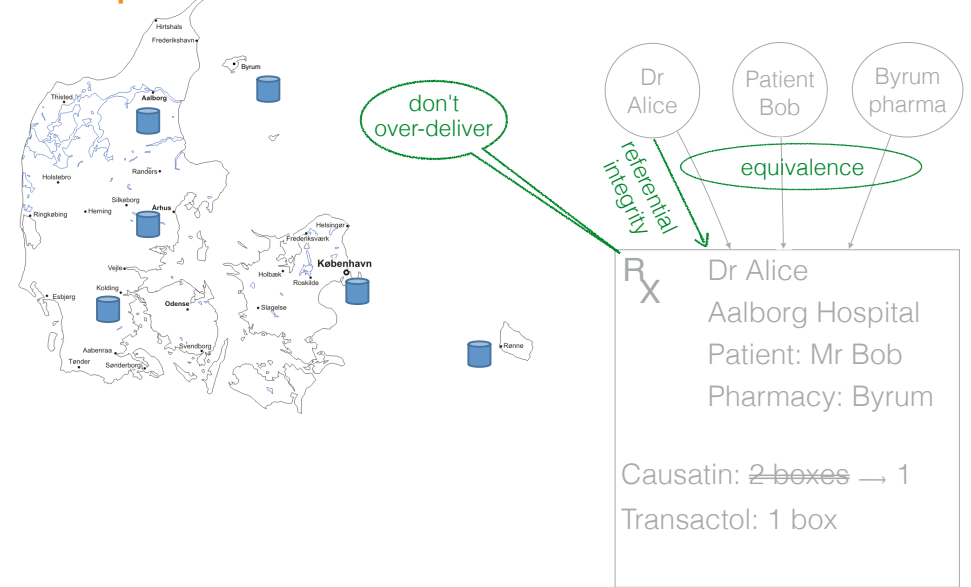# Classifying invariants by their coordination protocol
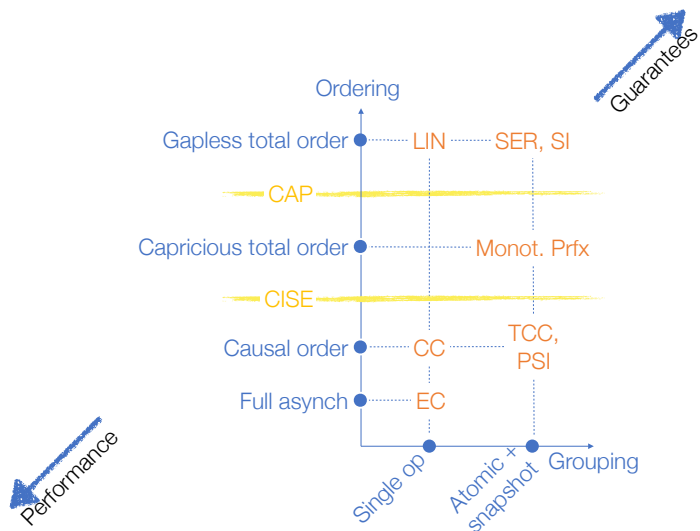
## FMKe *Fælles Medicinkort*



Dr Alice

Patient Bob

Byrum pharma

create (…)

add-med (…)

del-med (…)

R<sub>X</sub>

Dr Alice
Aalborg Hospital
Patient: Mr Bob
Pharmacy: Byrum

Causatin: 2 boxes → 1
Transactol: 1 box

# FMKe invariants



don't over-deliver

referential integrity

equivalence

create (…)

add-med (…)

del-med (…)

$$\text{add-med(n)} \quad \frac{r \geq 0 \qquad n \geq 0 \qquad r' \geq 0}{r' = r+n}$$

$$\text{del-med(n)} \quad \frac{r \geq 0 \qquad r \geq n \geq 0 \qquad r' \geq 0}{r' = r-n}$$

Rx

Dr Alice

Aalborg Hospital

Patient: Mr Bob

Pharmacy: Byrum

Causatin: ~~2 boxes~~ → 1

Transactol: 1 box

Dr Alice — Patient Bob — Byrum pharma

# Replicated FMKe: invariants?



don't over-deliver

referential integrity

equivalence

Rx

Dr Alice

Aalborg Hospital

Patient: Mr Bob

Pharmacy: Byrum

Causatin: ~~2 boxes~~ → 1

Transactol: 1 box

Dr Alice — Patient Bob — Byrum pharma

# What protocols for what invariants?



Ordering

Guarantees

Gapless total order — LIN — SER, SI

CAP

Capricious total order — Monot. Prfx

CISE

Causal order — CC — TCC, PSI

Full asynch — EC

Performance

Single op — Atomic + snapshot — Grouping

# Fully commutative updates

Some examples:
- Non-shared state
- Local blind: e.g., credit card number $hash(n) = 0$
- Empty invariant + fully-commuting CRDTs
  - LWW
  - Grow-only set
  - PN counter
  - Vector clock

Convergent

Asynchronous propagation
- Perfect scalability
- Perfect availability under partition

## What protocols for fully commutative?

## A ⟺ B: transactions



all-or-nothing

create (…)

R
 X    Dr Alice
      Aalborg Hospital
      Patient: Mr Bob
      Pharmacy: Byrum

## Transaction: Atomic writes + snapshot reads

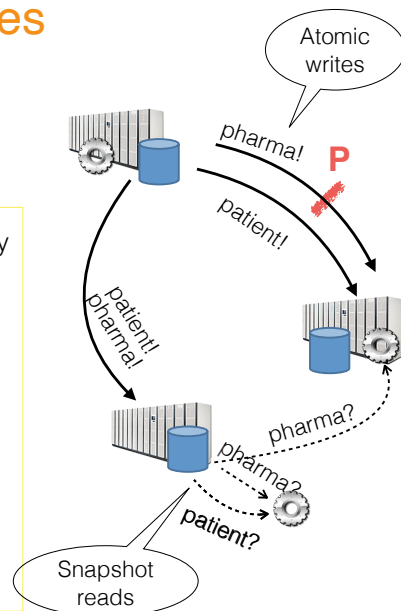*create-p* updates doctor, patient & pharmacy record

Atomic:
- = All-or-Nothing (A of ACID)
- Transmit joint updates together
- asynchronous

Snapshot: single database state
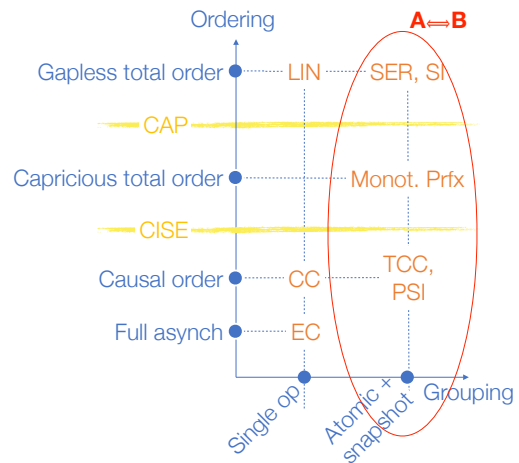- multi-version concurrency control
- asynchronous

Asynch: Available under partition

## "A ⟺ B" style invariants

Some example cases:
- *A=B*
- *A = ¬B*
- *friend(x,y) ⟺ friend(y,x)*
- *x≤y: <x++;y++>* (one actor)

# What protocols for A ⟸⟹ B?



Ordering

Gapless total order — LIN ⋯ SER, SI

A⟸⟹B

CAP

Capricious total order — Monot. Prfx

CISE

Causal order — CC — TCC, PSI

Full asynch — EC

Single op — Atomic + snapshot — Grouping

# A ⟹ B: demarcation



Dr Alice   Patient Bob   Byrum pharma

relative order

**2**

**1** create (…)

R X   Dr Alice
Aalborg Hospital
Patient: Mr Bob
Pharmacy: Byrum

# Distributed checkpoint: demarcation



checkpoint          hc

lj          journal
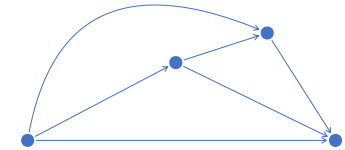
$Inv(\sigma) = \{ \sigma.lj \le \sigma.hc\}$       // no gap

Demarcation protocol

- Journal dæmon: lower bound += d
- Journal: Send message to Checkpoint *advance(d)*
- Checkpoint dæmon: upper bound += d

# Distributed checkpoint: demarcation



checkpoint          hc

*advance(d)*

lj          journal

$Inv(\sigma) = \{ \sigma.lj \le \sigma.hc\}$       // no gap

Demarcation protocol

- Journal dæmon: lower bound += d
- Journal: Send message to Checkpoint *advance(d)*
- Checkpoint dæmon: upper bound += d

# Distributed checkpoint: demarcation



$Inv(\sigma) = \{ \sigma.lj \leq \sigma.hc\}$     // no gap

Demarcation protocol
- Journal dæmon: lower bound += d
- Journal: Send message to Checkpoint *advance(d)*
- Checkpoint dæmon: upper bound += d

# Non-commuting, demarcation

Non-commuting CRDTs
- Empty invariant
  - Set
  - Map
  - MVR
- Universally-stable operations
  - Acyclic graph: add-parallel, remove
  - Sequence: insert-at, remove

2-actor implication invariants
- $A \Longrightarrow B$
- Referential integrity
- Chicken/fox/grain: *grain* $\Longrightarrow$ *¬chicken*
- $x \leq y$

Preserve order across processes
  → Causal Consistency

# Causal Consistency

*create-p* before *add-med*
- "Bob points to Rx $\Longrightarrow$ Rx valid"
  - Referential integrity
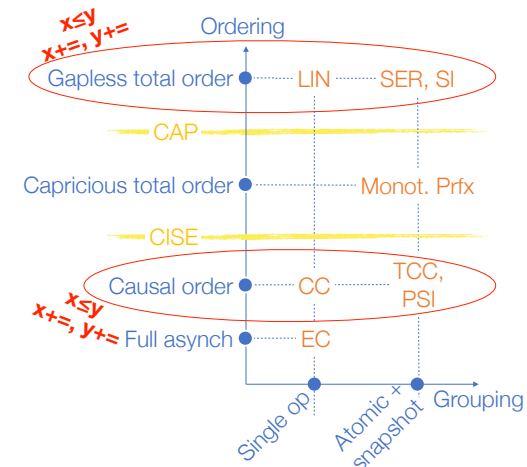- General case: LHS $\Longrightarrow$ RHS
- pattern: RHS!; LHS!

Deliver in the right order: Causal Consistency

Local decision:
- requires metadata
- available

# What protocols for demarcation?

## x≥0 x—: total/mutual order

*u(), v()* not mutually stable
- "Conflicting"
- Either *u()* before *v()*, or *v()* before *u()*

Protocols:
- General case: total order, consensus
- 1 lock / set of mutually-conflicting operations
  - Coarser locks OK
- Single server / conflict set (flat combining)
- Social

---

## Bounded Counter

Specific, common case

Shared counter:
- $x \geq 0$
- *increment (n)*
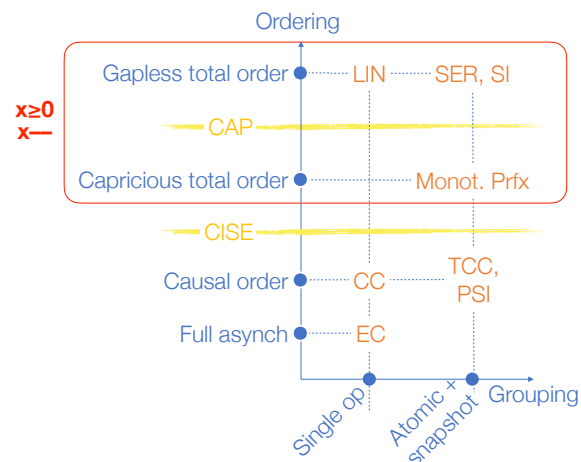- *decrement (n)*                    // precondition $x \geq n$

Escrow:
- Local share, decrement *share −= n*
- decrement disallowed if *share < n*
- Donate share

Mostly AP

Encapsulated, proven correct (CISE)

Causal ordering essential

---

## What protocols for x≥0 x—?

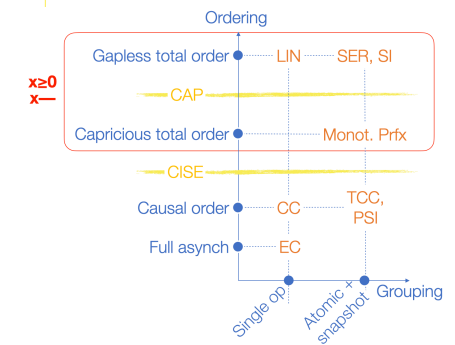---

## Mutual order: *a posteriori* vs. *a priori*

*A posteriori,* Capricious, Monotonic-Prefix:
- Execute
- Pick a number
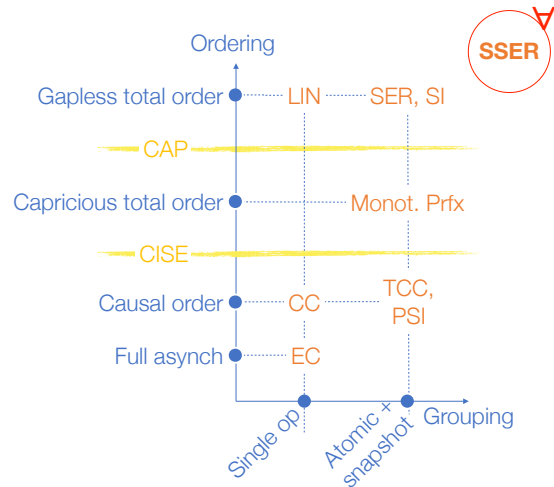- Propagate
- Sort
- Roll back; roll forward
- Iterate

*A priori*, Gapless, Monotonic Reads:
- Consensus on a number
- Wait for my turn
- Execute & propagate

Capricious + finality ⟹ consensus

# What protocols for arbitrary invariants?

Ordering

SSER ∀

- Gapless total order — LIN — SER, SI
- CAP
- Capricious total order — Monot. Prfx
- CISE
- Causal order — CC — TCC, PSI
- Full asynch — EC

Single op · Atomic + snapshot · Grouping

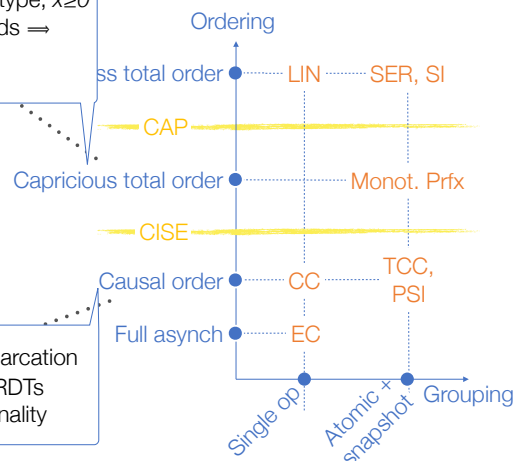# Sweet spot: Transactional Causal Consistency + optional consensus

TCC =
- Causal consistency
  - $x \leq y$, demarcation
- Snapshot reads + Atomic writes
  - $A \Longleftrightarrow B$
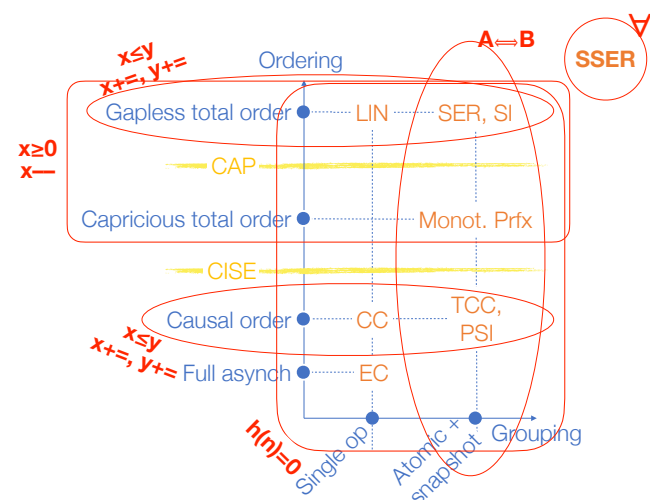- System: AntidoteDB

Available: not $x \geq 0$
- Strengthen when necessary
- System: Colony

# Strongest AP protocol(s)

Total order
Any sequential data type, $x \geq 0$
Non Monotonic Reads $\Rightarrow$ rollbacks
No finality

Ordering

- ss total order — LIN — SER, SI
- CAP
- Capricious total order — Monot. Prfx
- CISE
- Causal order — CC — TCC, PSI
- Full asynch — EC

Partial order $\longrightarrow$ demarcation
Requires Merge / CRDTs
Monotonic Reads; finality

Single op · Atomic + snapshot · Grouping

# What protocols for what invariants?

$x \leq y$
$x+=, y+=$
$A \Longleftrightarrow B$
SSER ∀

Ordering

- Gapless total order — LIN — SER, SI
- $x \geq 0$ $x$—
- CAP
- Capricious total order — Monot. Prfx
- CISE
- Causal order — CC — TCC, PSI
- $x \leq y$ $x+=, y+=$ Full asynch — EC
- h(n)=0

Single op · Atomic + snapshot · Grouping