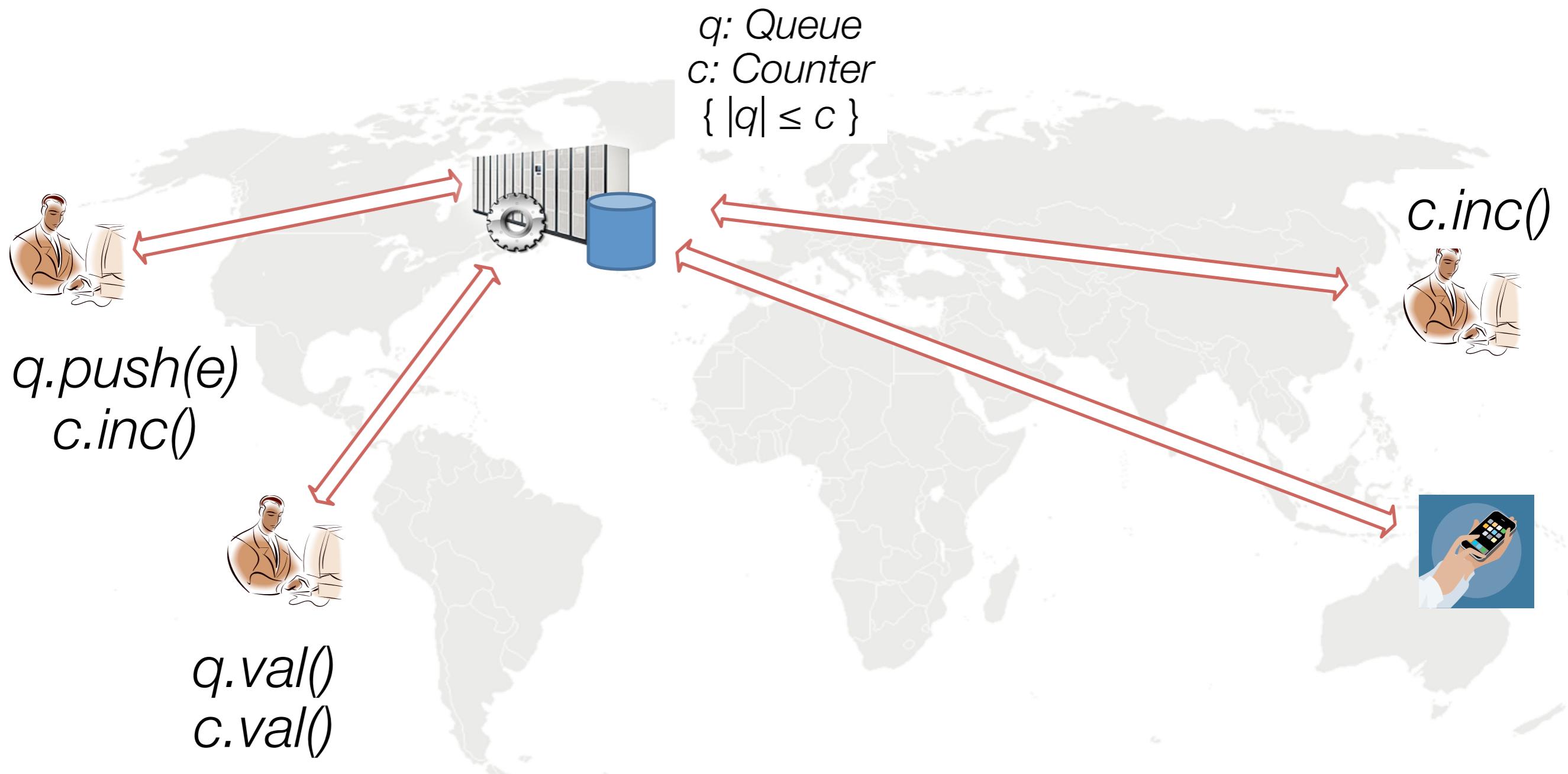
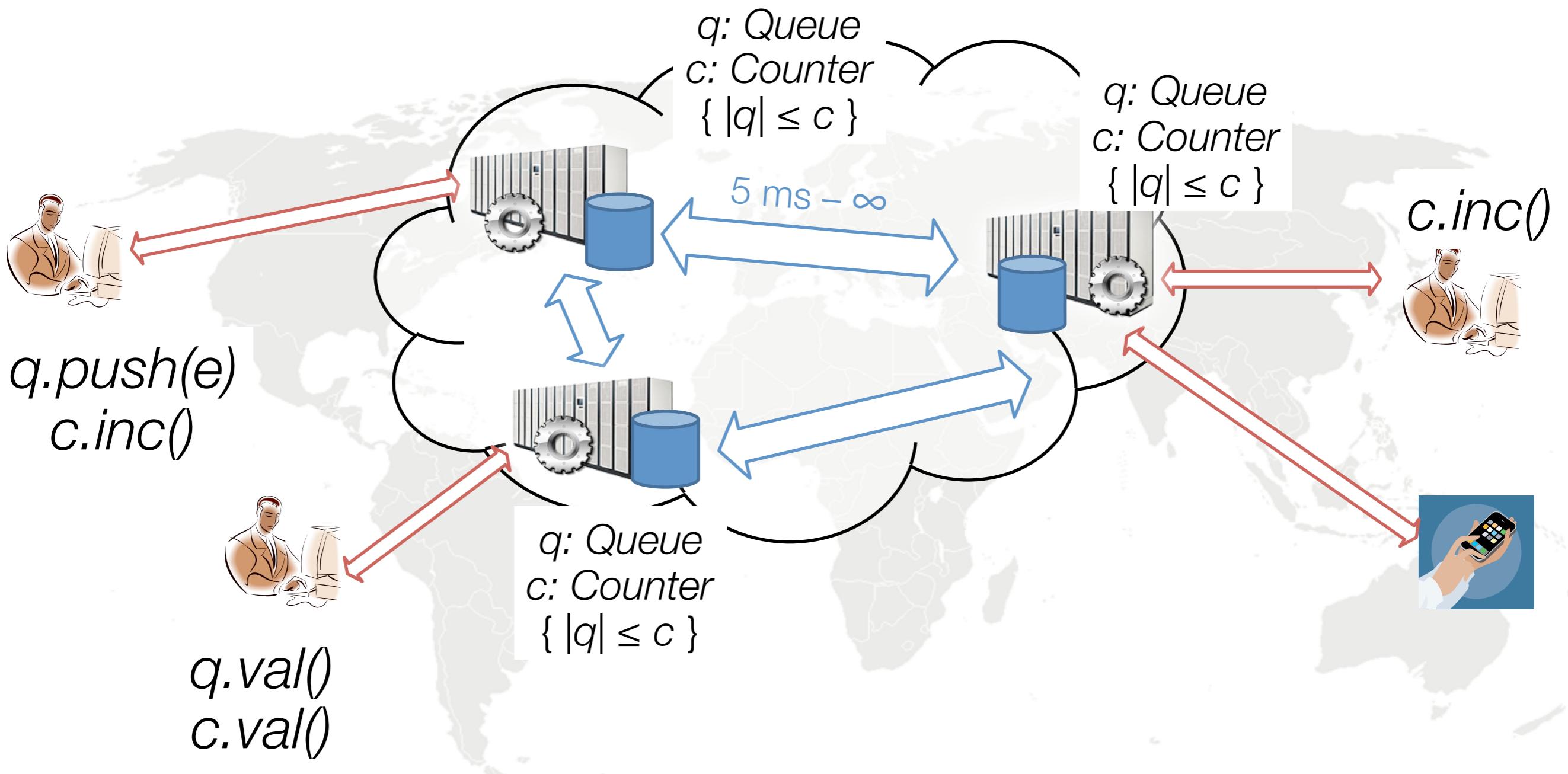


Distributed Consistency

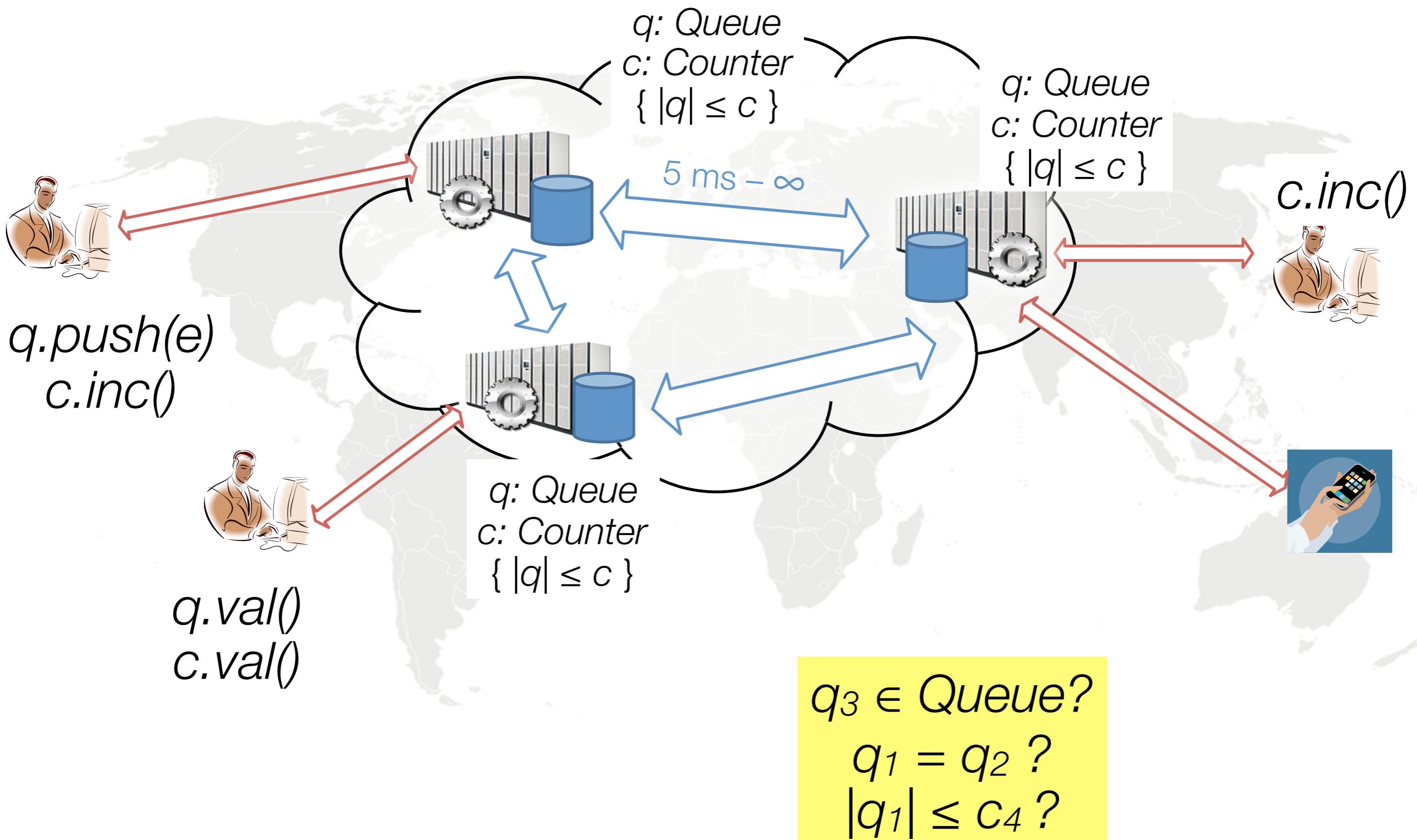
Shared database



Geo-replicated



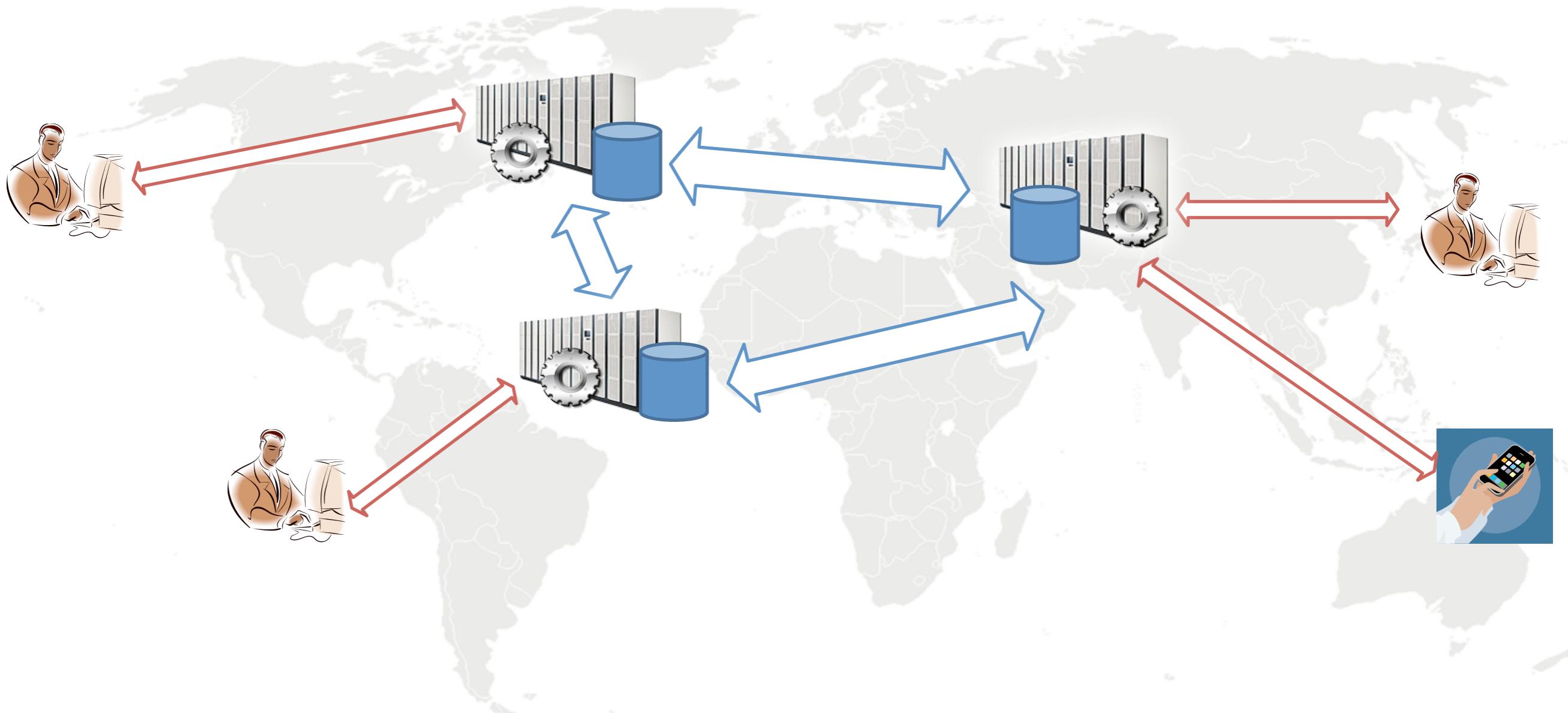
Geo-replicated



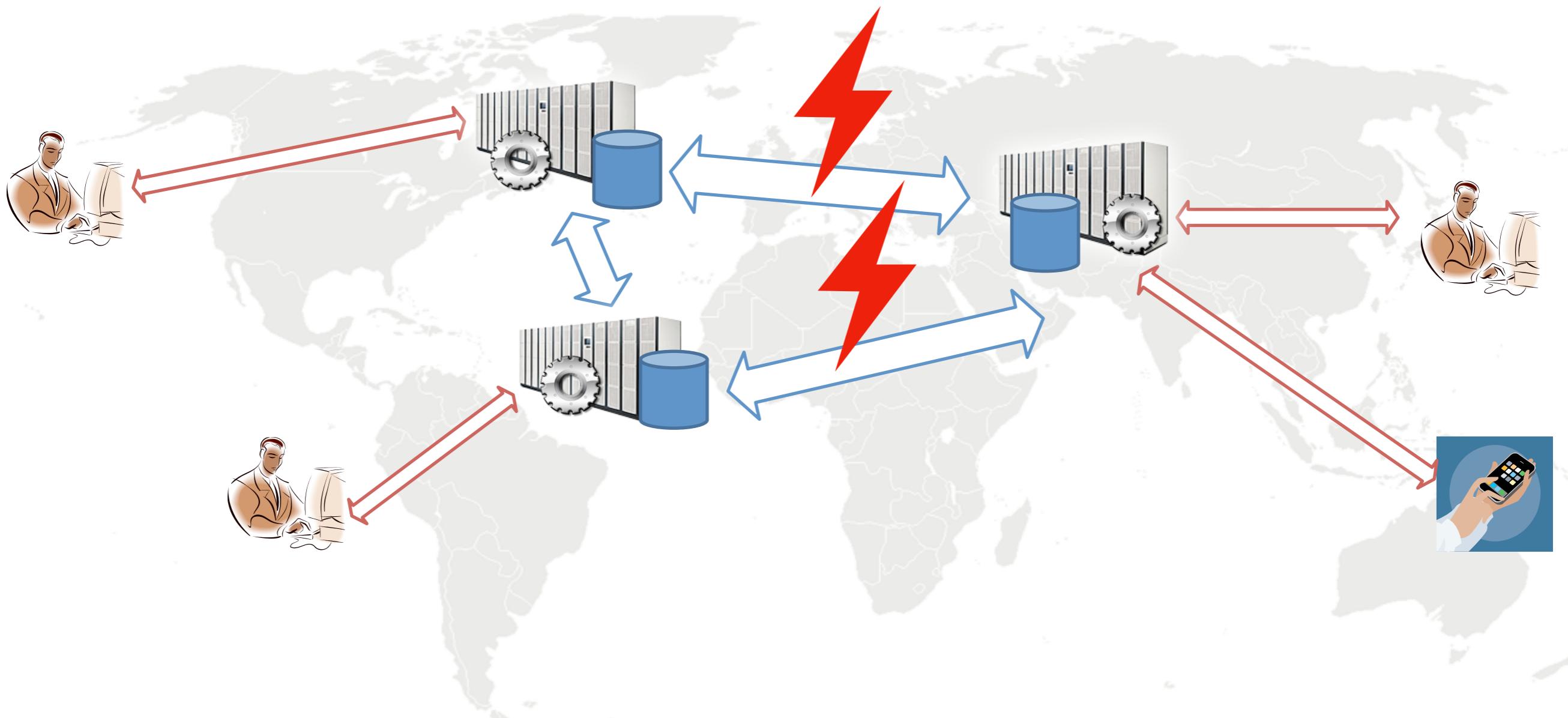
Consistency

- More replicas:
 - Better read availability, responsiveness, performance
 - More work to keep replicas in sync
- Consistent = behavior similar to sequential:
 - Satisfies specs: does q behave like a queue?
 - Replicas agree: is q identical everywhere?
 - Objects agree: is $|q| \leq c$?
 - Same flow of time? $q1.push()$ before $q2.push()$

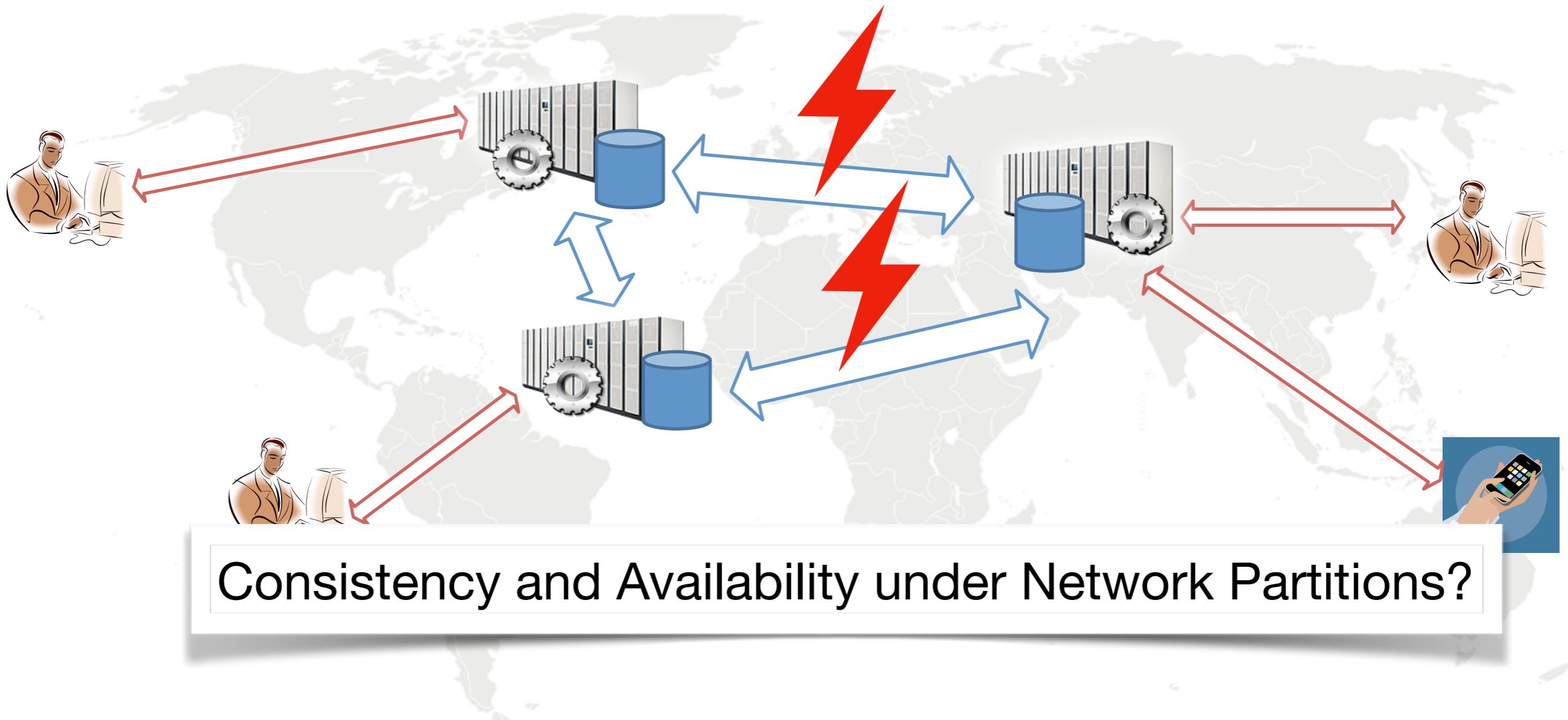
CAP Theorem



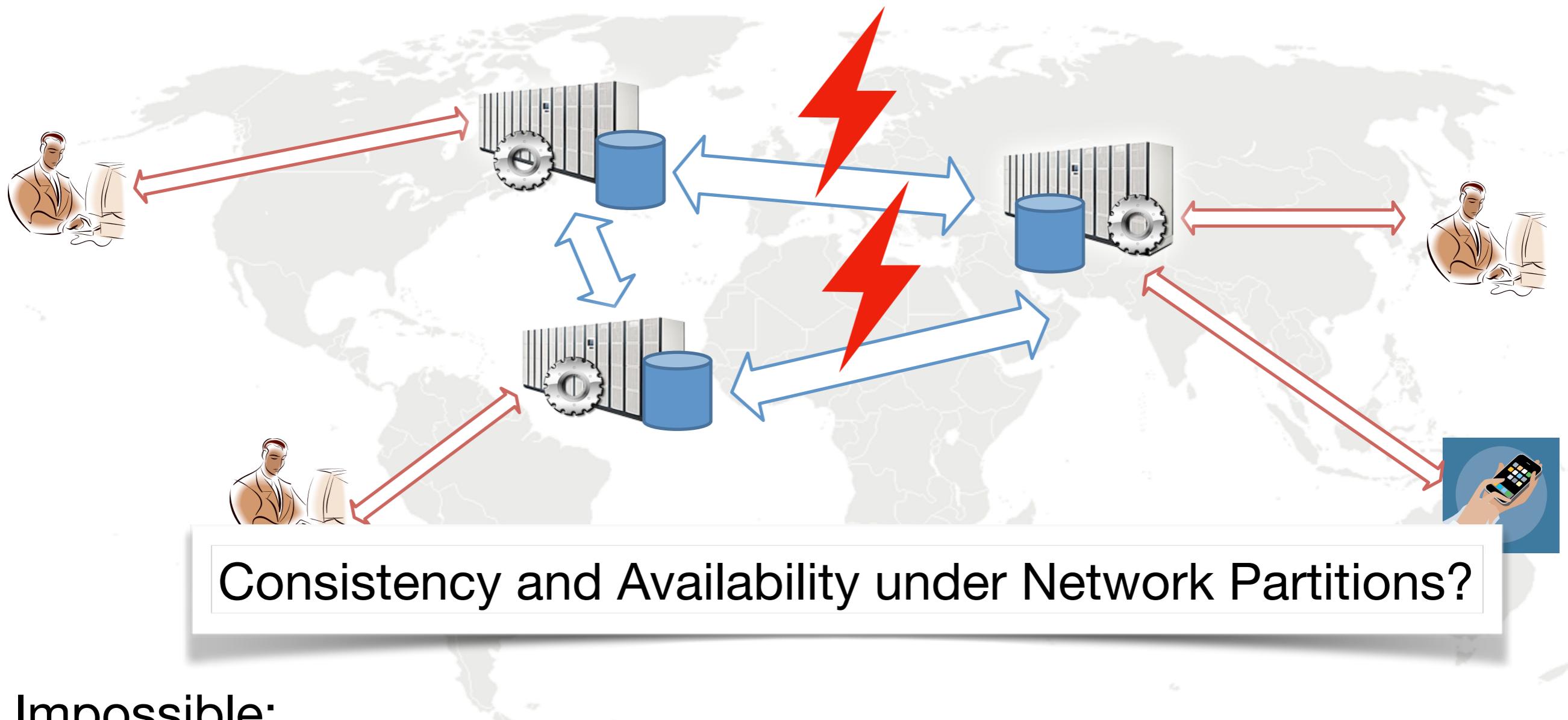
CAP Theorem



CAP Theorem



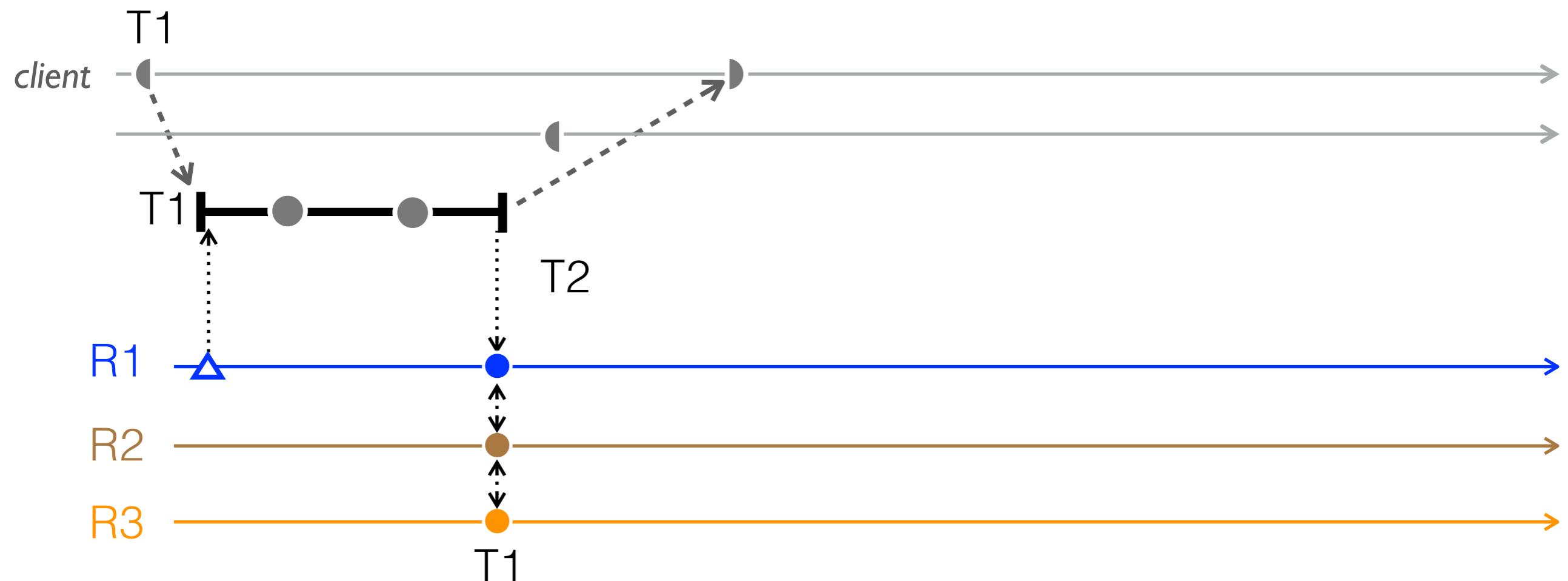
CAP Theorem



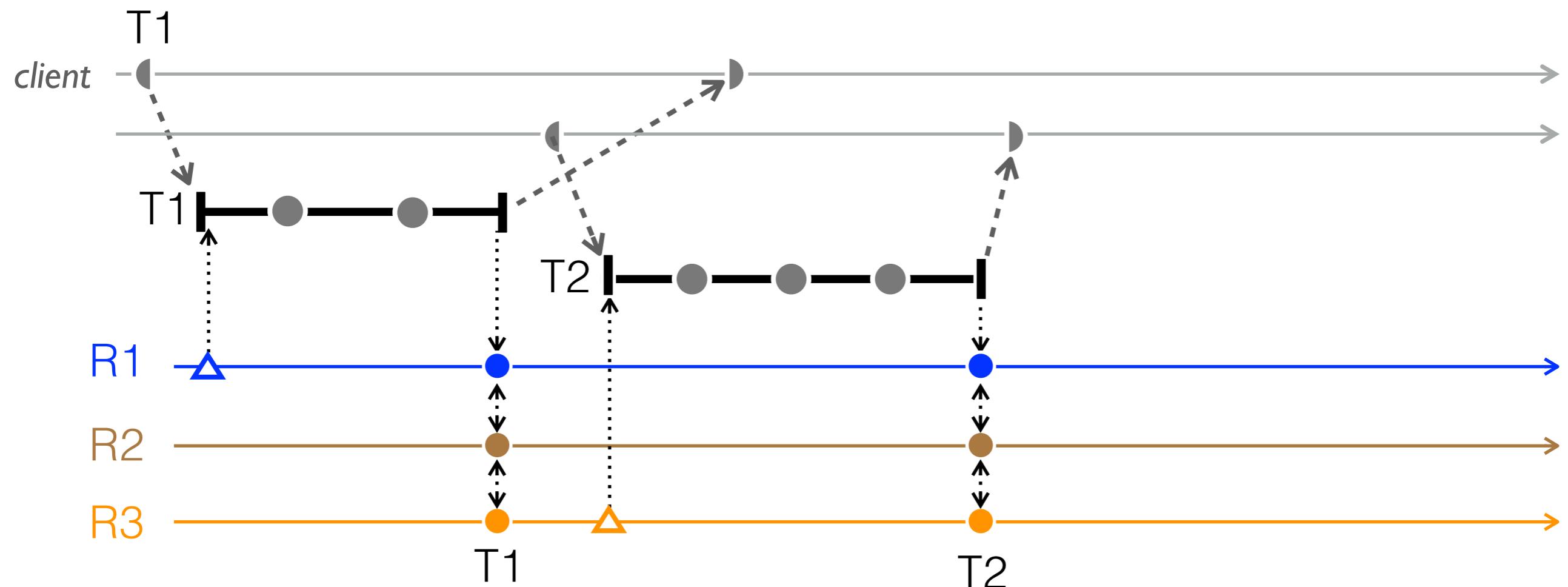
Impossible:

- Consistency: the system has to stop until the network is restored
- Availability: we have to let different replicas diverge (for a while)

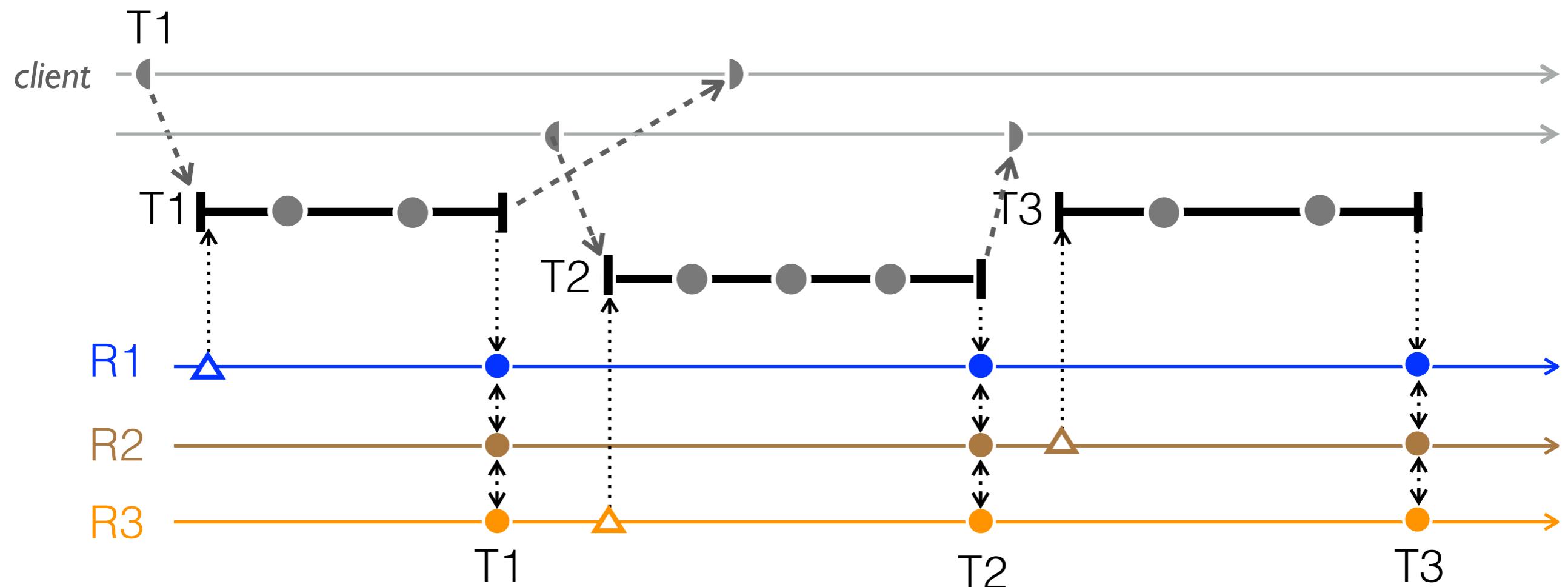
Strict Serialisability



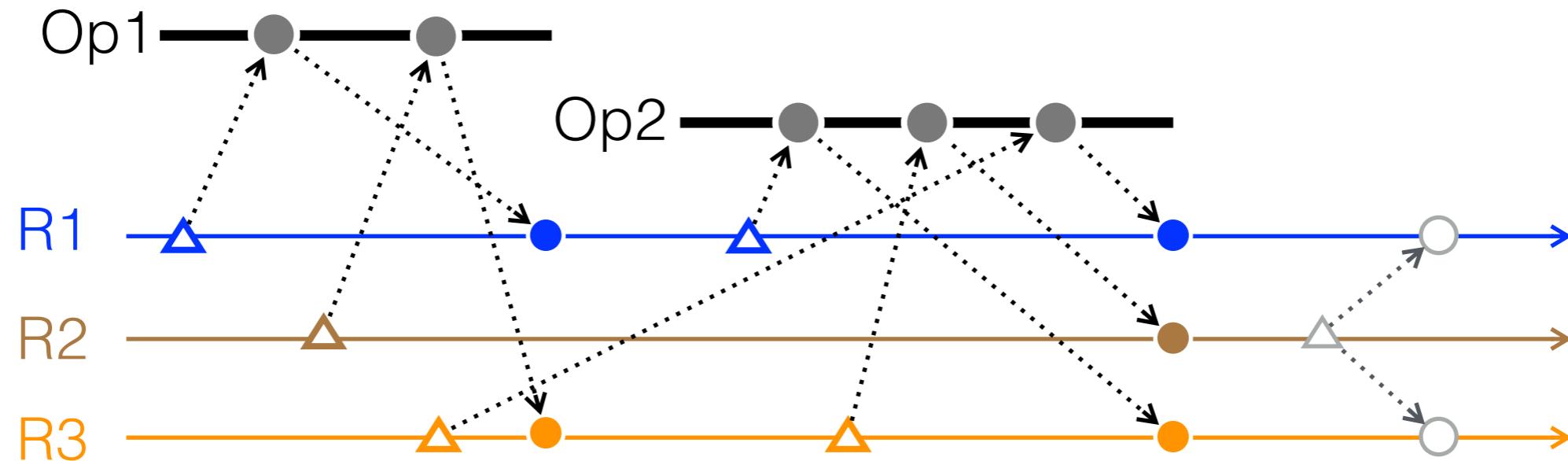
Strict Serialisability



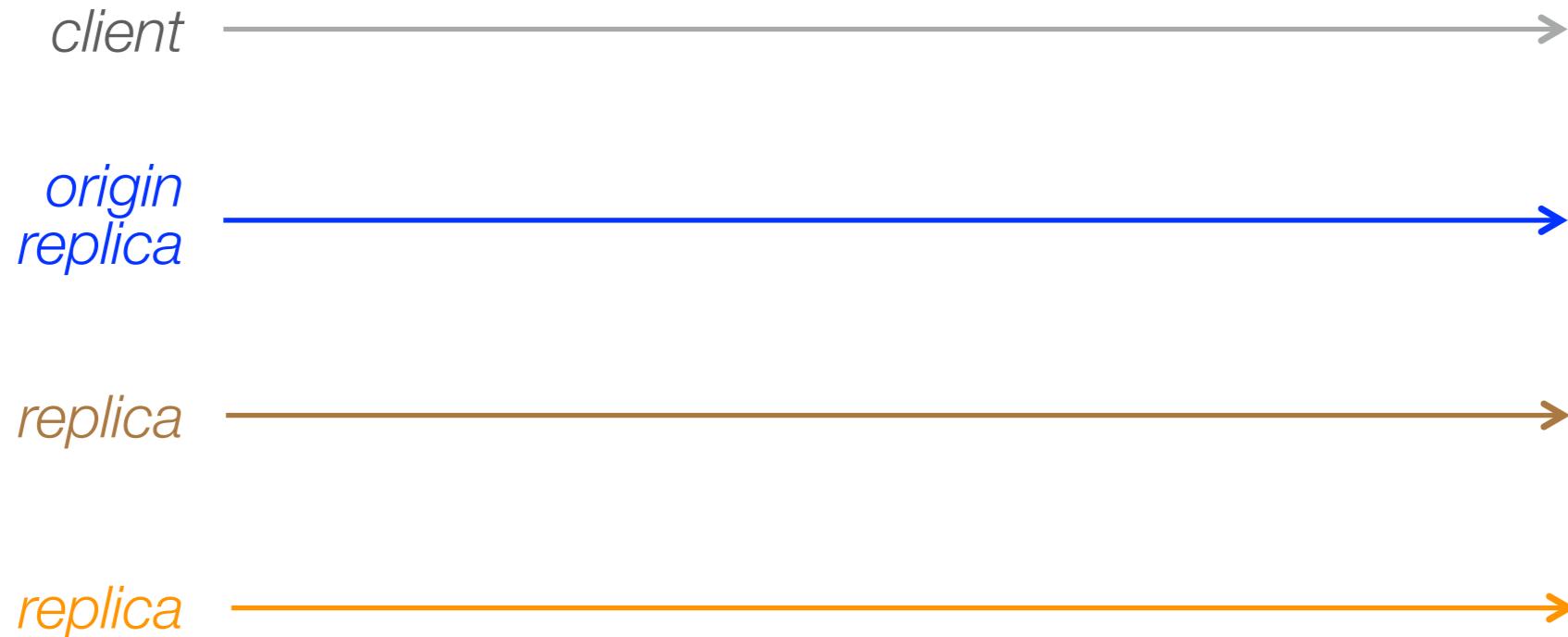
Strict Serialisability



Eventual consistency



Replicated operation



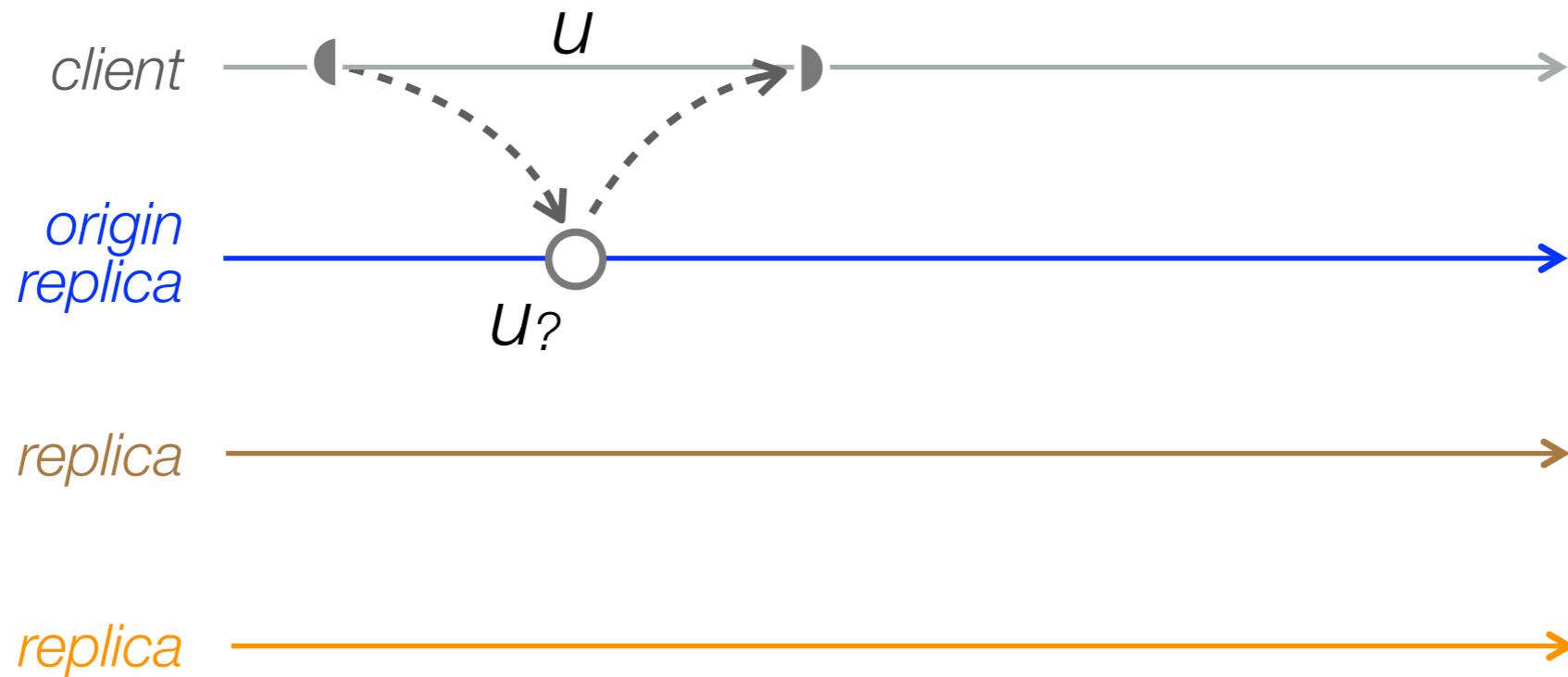
$u : state \hookrightarrow (retval, (state \hookrightarrow state))$

Prepare (@origin) $u?$; deliver $u!$

Read one, write all (ROWA)

Deferred-update replication (DUR)

Replicated operation



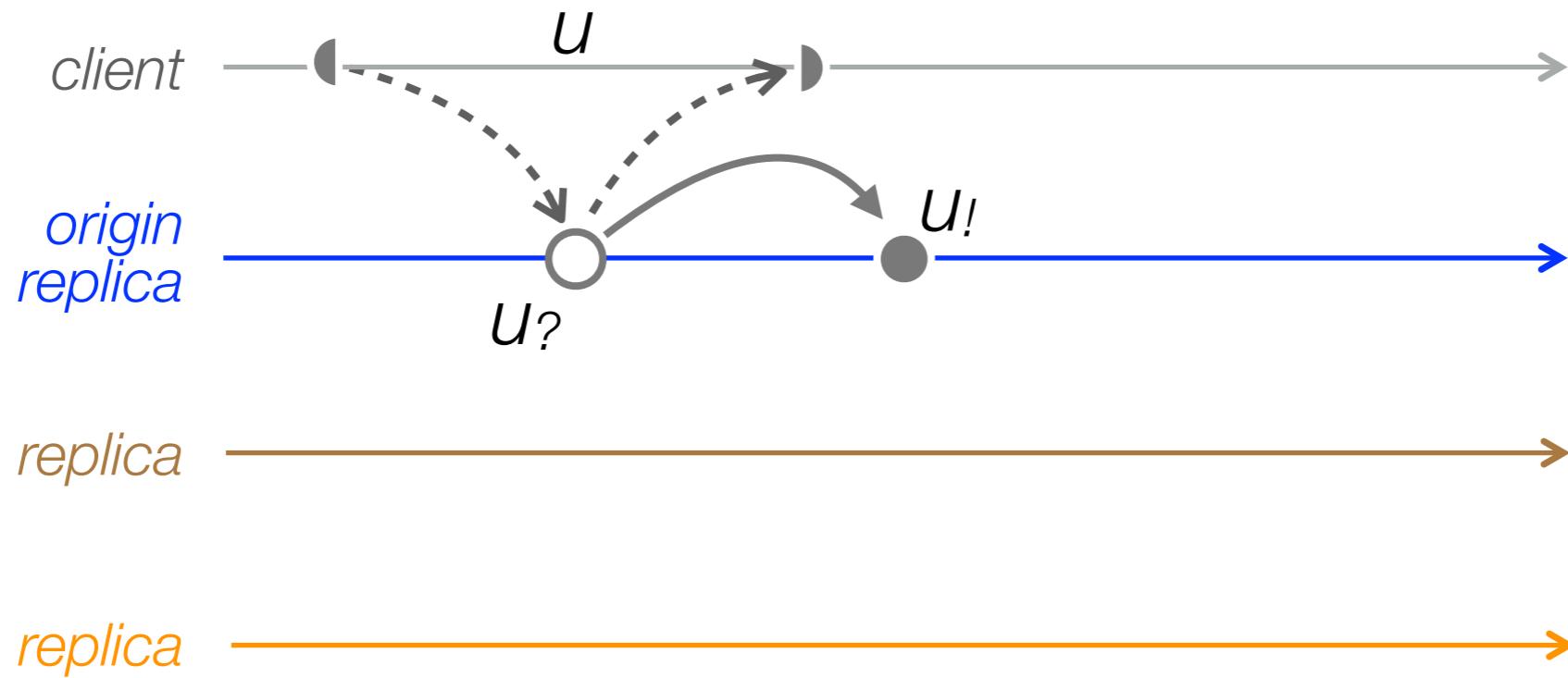
$u: \text{state} \rightsquigarrow (\text{retval}, (\text{state} \rightsquigarrow \text{state}))$

Prepare (@origin) $u?$; deliver $u!$

Read one, write all (ROWA)

Deferred-update replication (DUR)

Replicated operation



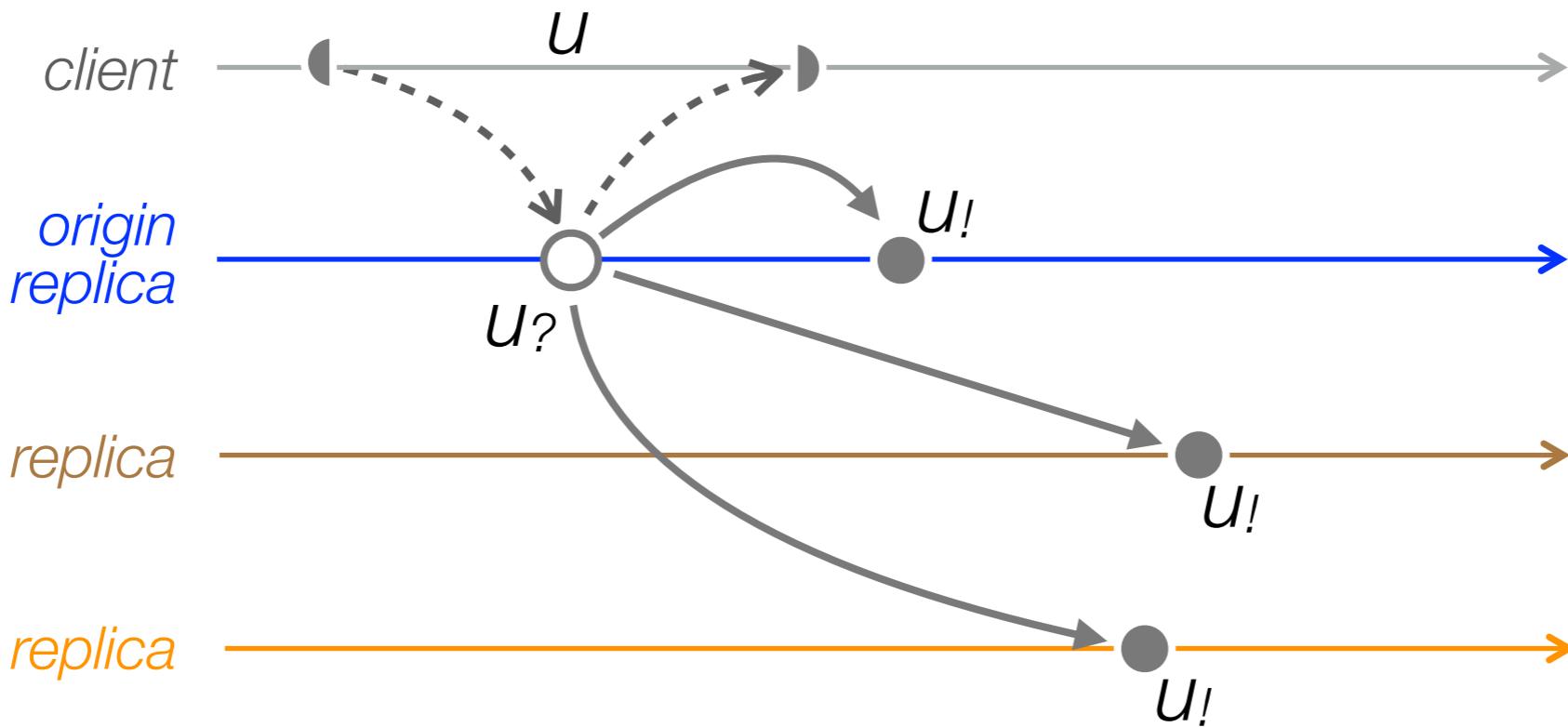
$u: \text{state} \rightsquigarrow (\text{retval}, (\text{state} \rightsquigarrow \text{state}))$

Prepare (@origin) $u?$; deliver $u!$

Read one, write all (ROWA)

Deferred-update replication (DUR)

Replicated operation



$u: \text{state} \rightsquigarrow (\text{retval}, (\text{state} \rightsquigarrow \text{state}))$

Prepare (@origin) $u?$; deliver $u!$

Read one, write all (ROWA)

Deferred-update replication (DUR)

In this Class

- We will adopt a client view of consistency
 - Just like with Memory Models
- Different protocols implement different consistency criteria
 - Stronger protocols are more expensive in performance but limit the non-determinism for clients
- We will use a uniform semantics for clients that over-approximates the non-determinism
 - [Burckhardt,Gotsman,Yang'15]

In this Class

- We will adopt a client view of consistency
 - Just like with Memory Models
- Different protocols implement different levels of consistency
 - Stronger protocols are more efficient but limit the non-determinism
- We will use a uniform semantics that approximates the non-determinism [Burckhardt et al.]

Foundations and Trends® in Programming Languages
Vol. 1, No. 1-2 (2014) 1–150
© 2014 S. Burckhardt
DOI: 10.1561/2500000011



Principles of
Eventual Consistency

Sebastian Burckhardt
Microsoft Research
sburkha@microsoft.com

Client View of the system

- We will explain in an axiomatic way (à la CAT) the possible results of each operation
- As in CAT, we will posit the existence of certain orders that explain why a behavior is possible
- We will describe Distributed Data Structure specifications by exploiting these orders
- Implementations of distributed data structures can be verified against these specifications
 - We will not talk about verification

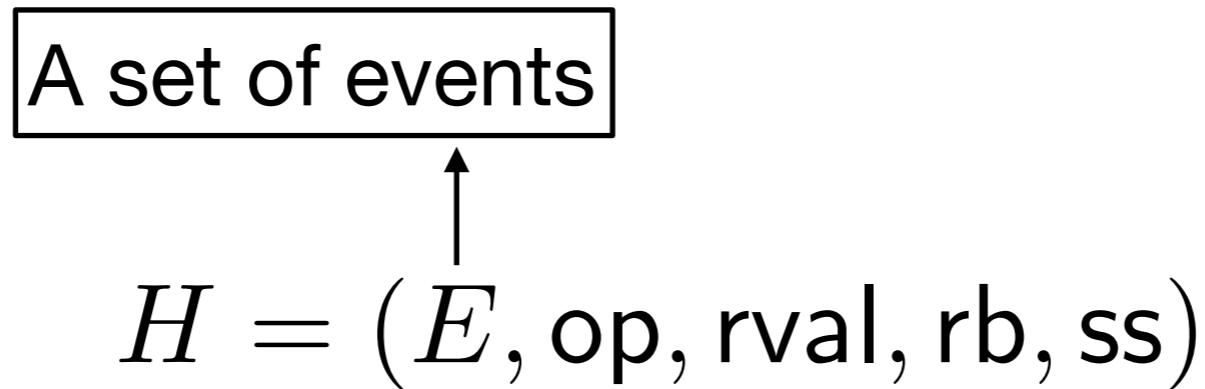
Client Operations

- Client submit operations which can in turn be transactions
- A client is represented as a Session
- A single session could issue multiple operations and transactions
- We will consider a session to be the equivalent of the program order in relaxed memory models

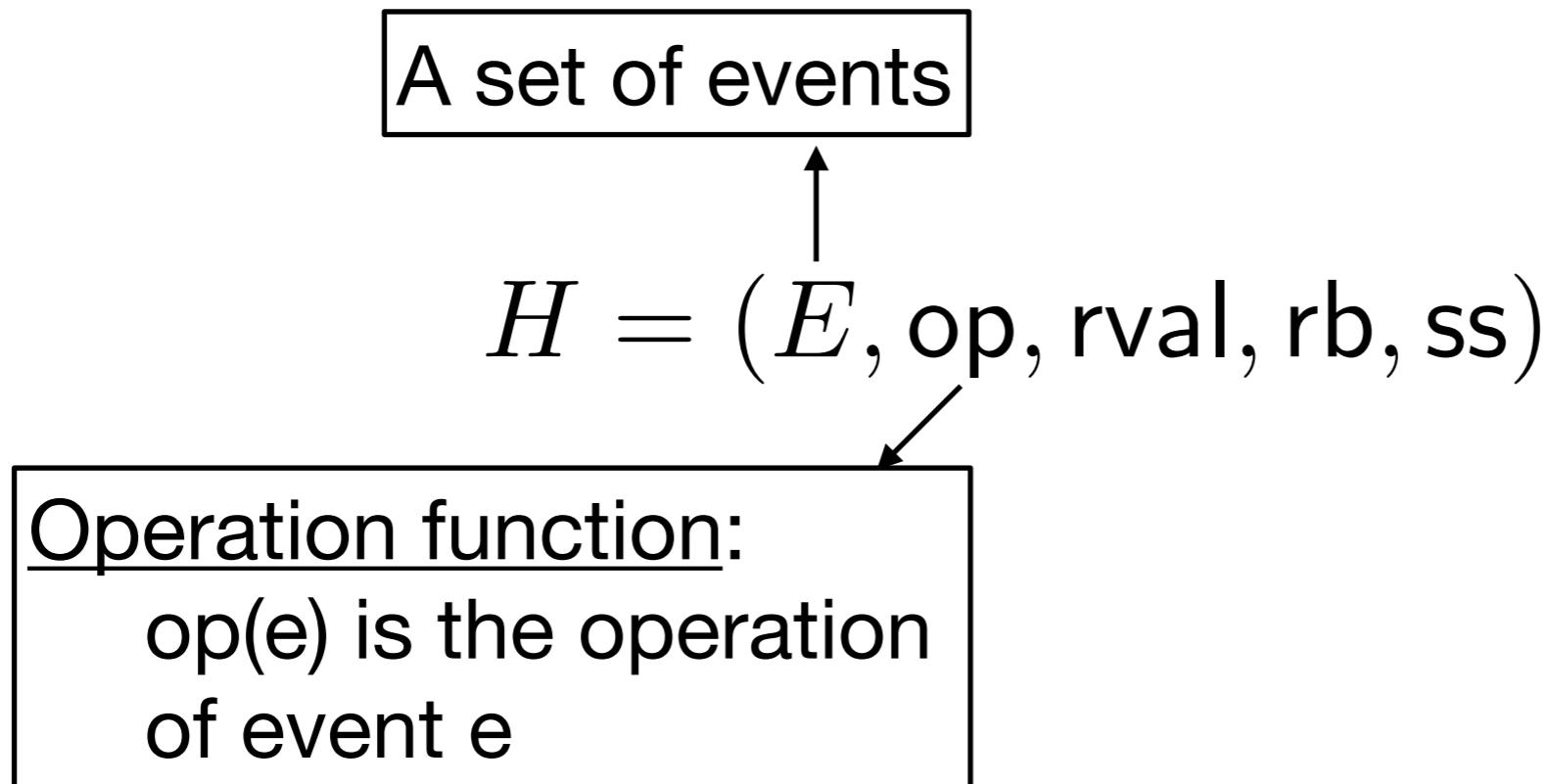
Axiomatic Consistency Histories

$$H = (E, \text{op}, \text{rval}, \text{rb}, \text{ss})$$

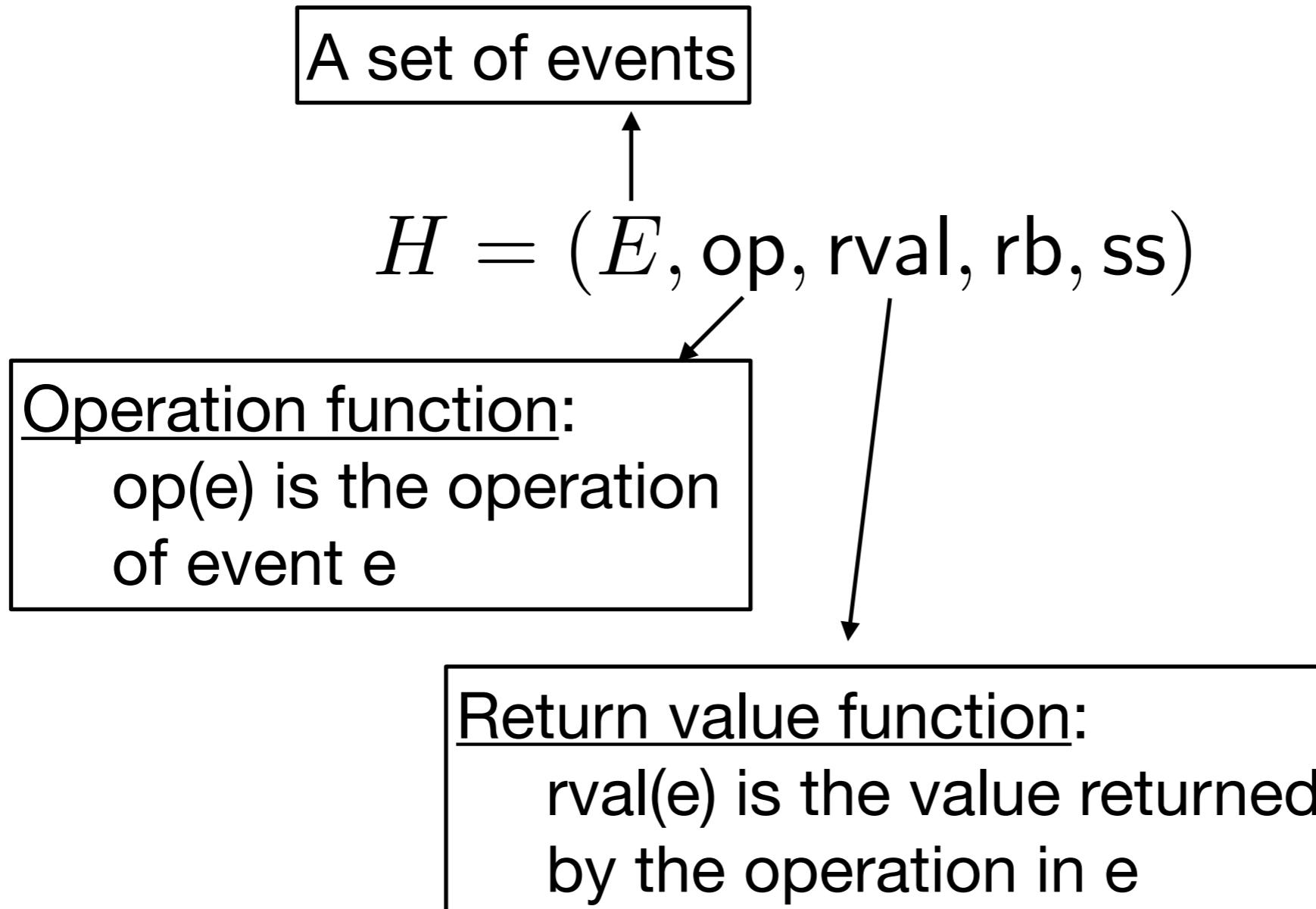
Axiomatic Consistency Histories



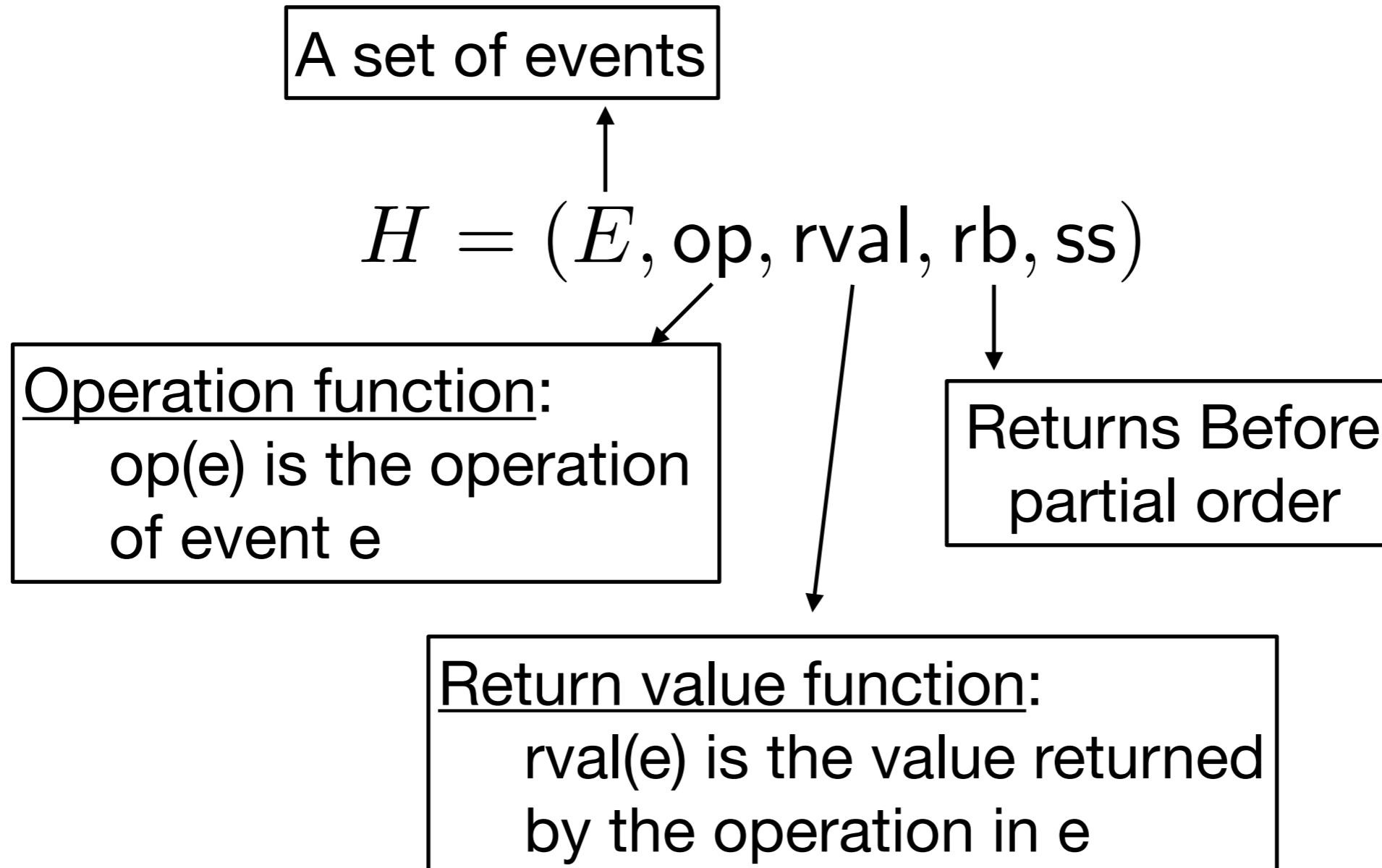
Axiomatic Consistency Histories



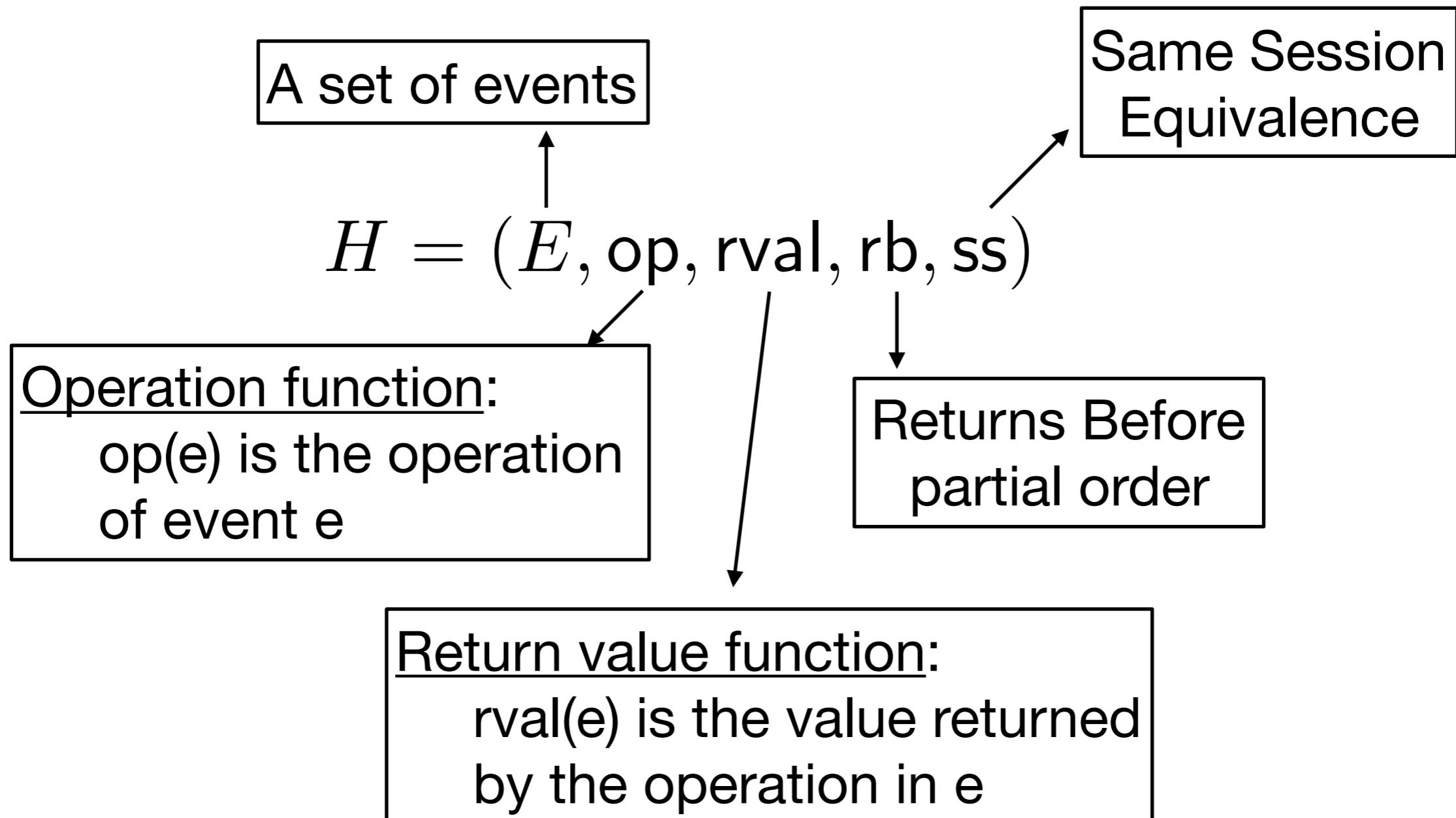
Axiomatic Consistency Histories



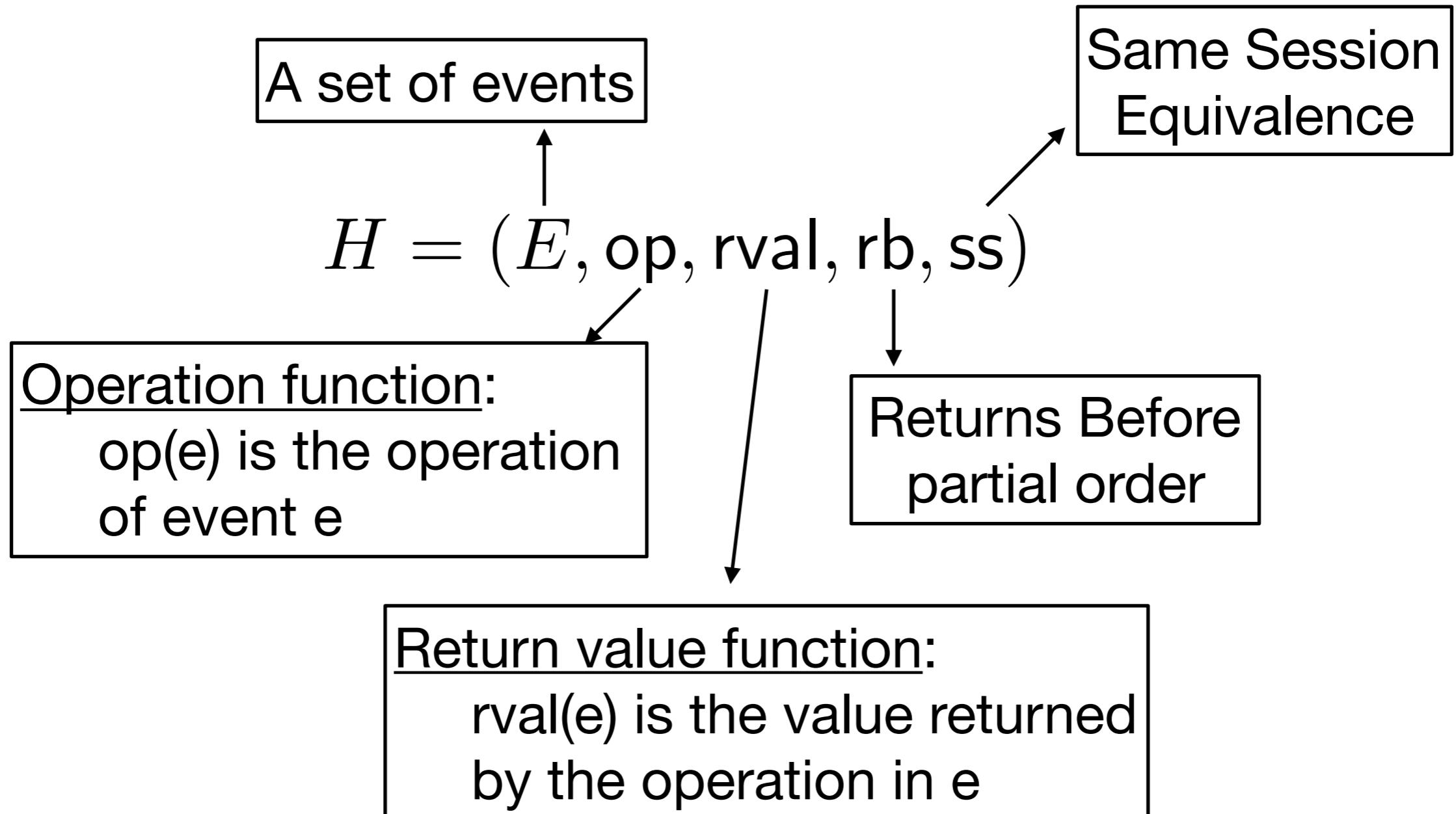
Axiomatic Consistency Histories



Axiomatic Consistency Histories

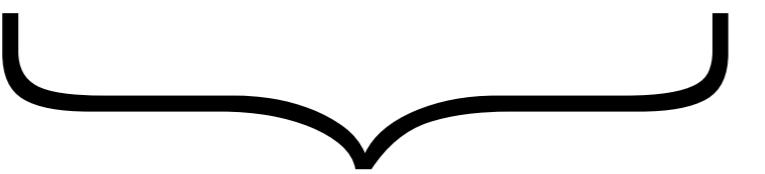


Axiomatic Consistency Histories



Session order $so = ss \cap rb$

Axiomatic Consistency Abstract Executions

$$A = (E, \text{op}, \text{rval}, \text{rb}, \text{ss}, \text{vis}, \text{ar})$$


History

Axiomatic Consistency Abstract Executions

Visibility partial order:
 $e \xrightarrow{\text{vis}} e'$ represents that
the effects of e are available
when executing e'

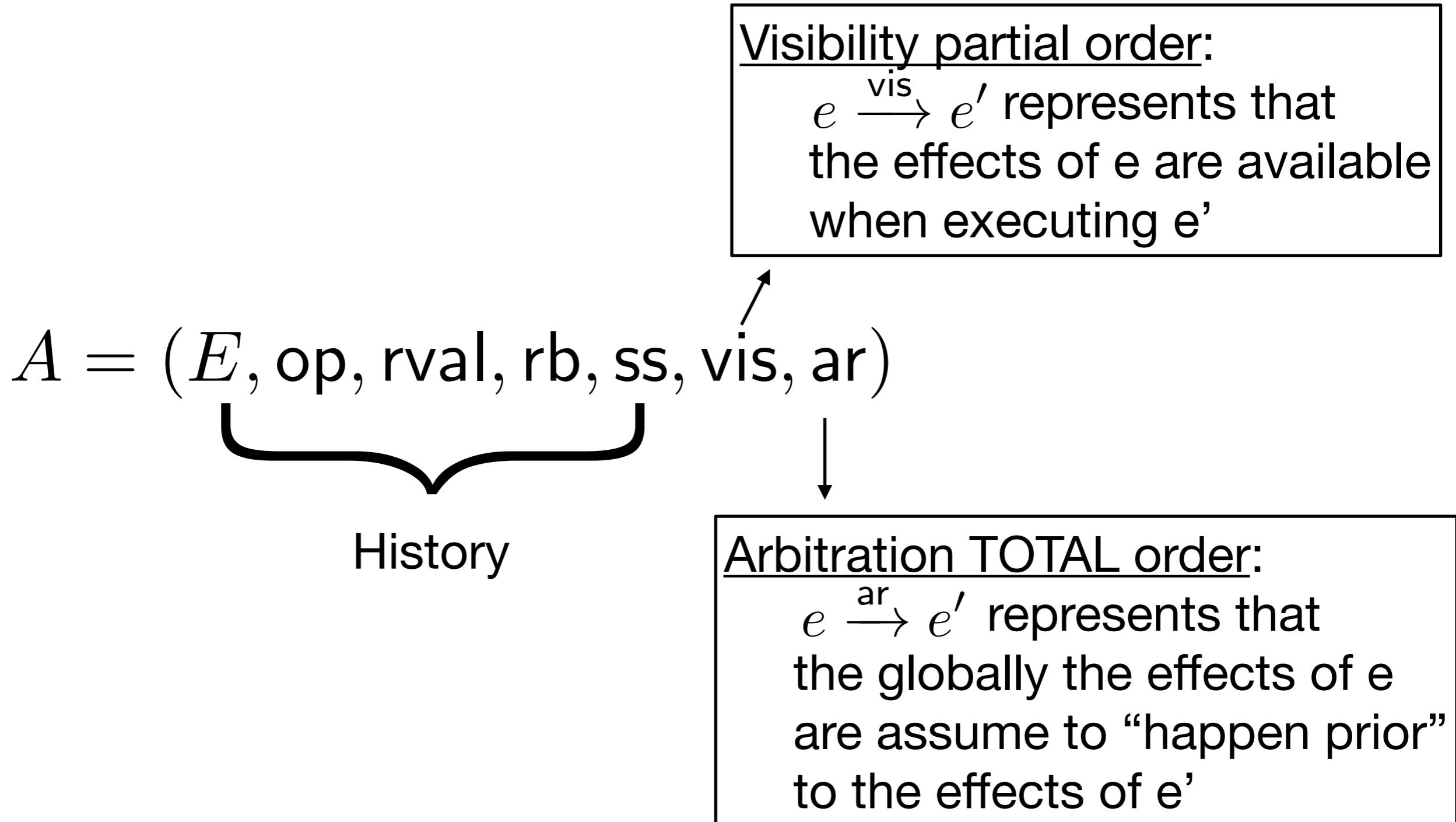
$$A = (E, \text{op}, \text{rval}, \text{rb}, \text{ss}, \text{vis}, \text{ar})$$

History



Axiomatic Consistency

Abstract Executions



Session Guarantees (Anomalies)

Session Guarantees for Weakly Consistent Replicated Data

Douglas B. Terry, Alan J. Demers, Karin Petersen, Mike J. Spreitzer, Marvin M. Theimer,
and Brent B. Welch

Computer Science Laboratory
Xerox Palo Alto Research Center
Palo Alto, California 94304

Abstract

Four per-session guarantees are proposed to aid users and applications of weakly consistent replicated data: Read Your Writes, Monotonic Reads, Writes Follow Reads, and Monotonic Writes. The intent is to present individual applications with a view of the database that is consistent with their own actions, even if they read and write from various, potentially inconsistent servers. The guarantees can be layered on existing systems that employ a read-any/write-any replication scheme while retaining the principal benefits of such a scheme, namely high-availability, simplicity, scalability, and support for disconnected operation. These session guarantees were developed in the context of the Bayou project at Xerox PARC in which we are designing and building a replicated storage system to support the needs of mobile computing users who may be only intermittently connected.

1. Introduction

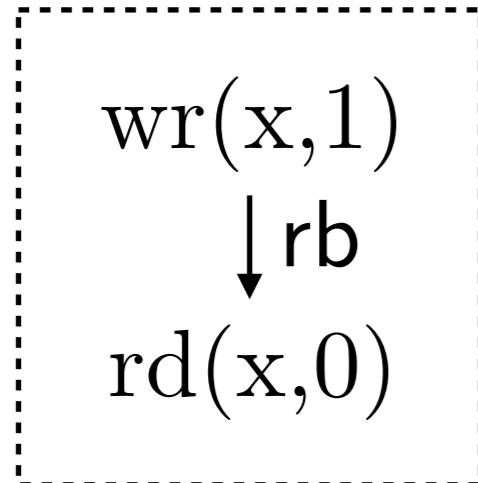
Techniques for managing weakly consistent replicated data have been employed in a variety of systems characterized by the

may want to read and update data copied onto their portable computers even if they did not have the foresight to lock it before either a voluntary or an involuntary disconnection occurred. Also, the presence of slow or expensive communications links in the system can make maintaining closely synchronized copies of data difficult or uneconomical.

Unfortunately, the lack of guarantees concerning the ordering of read and write operations in weakly consistent systems can confuse users and applications, as reported in experiences with Grapevine [21]. A user may read some value for a data item and then later read an older value. Similarly, a user may update some data item based on reading some other data, while others read the updated item without seeing the data on which it is based. A serious problem with weakly consistent systems is that inconsistencies can appear even when only a single user or application is making data modifications. For example, a mobile client of a distributed database system could issue a write at one server, and later issue a read at a different server. The client would see inconsistent results unless the two servers had synchronized with one another sometime between the two operations.

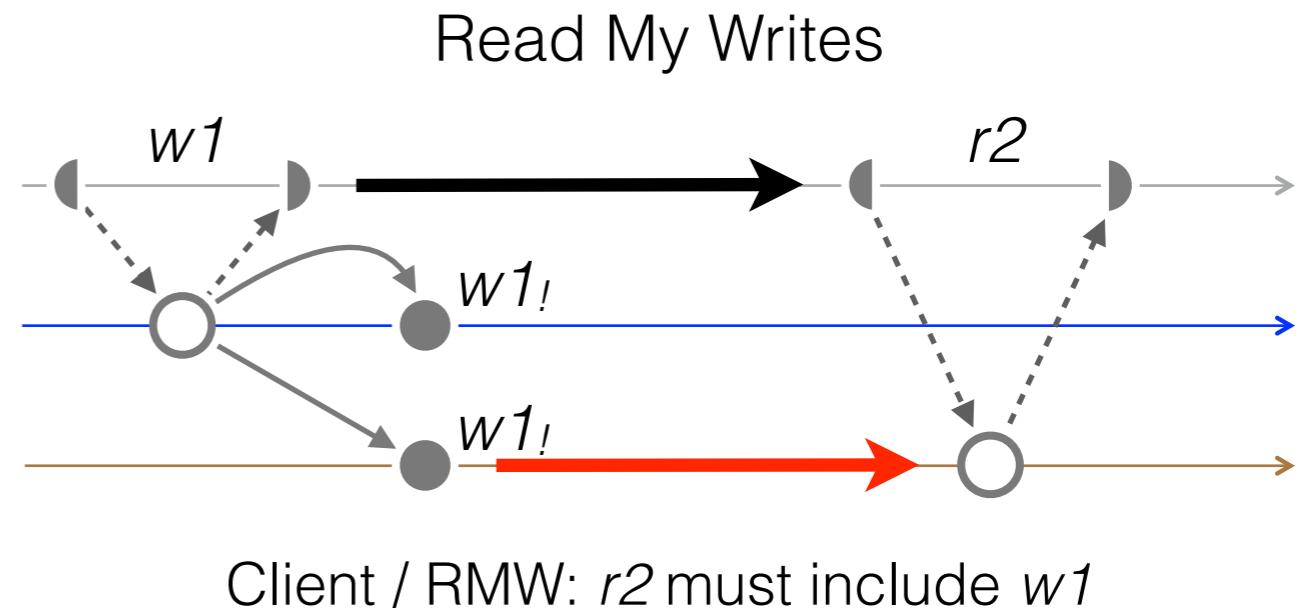
In this paper, we introduce *session guarantees* that alleviate this problem of weakly consistent systems while maintaining the principle advantages of read-any/write-

Session Guarantees (Anomalies)

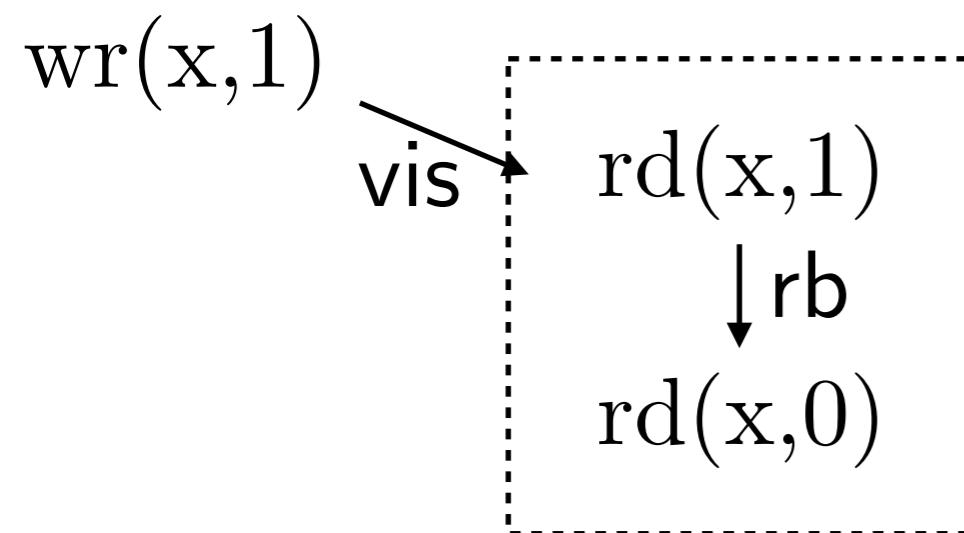
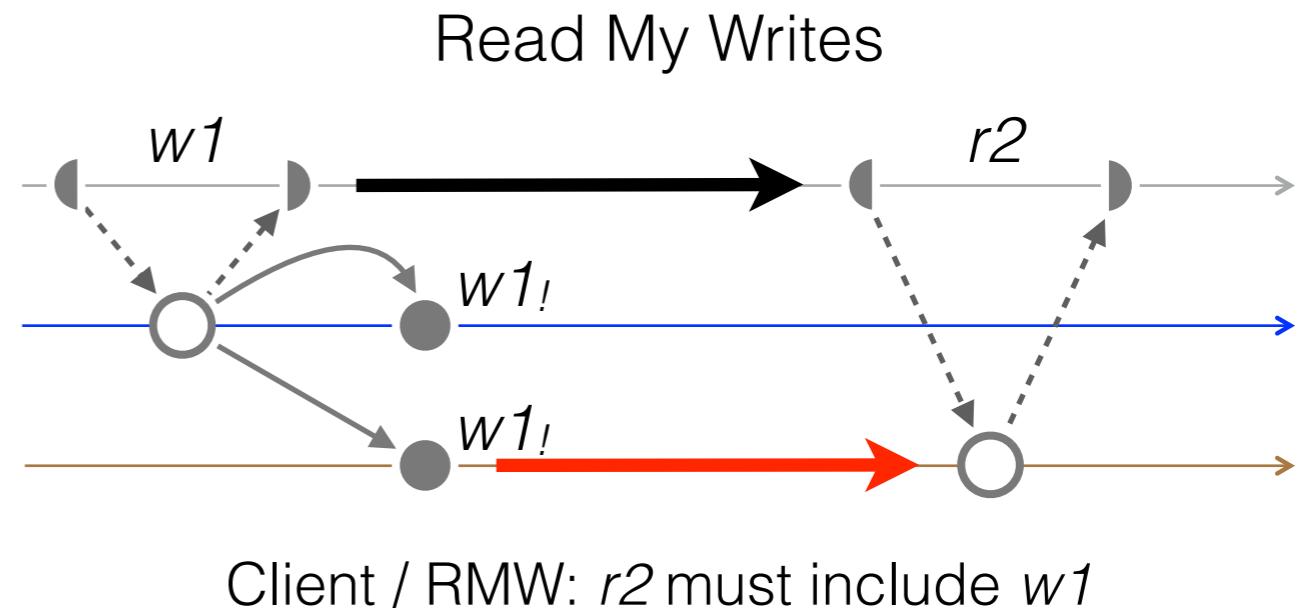
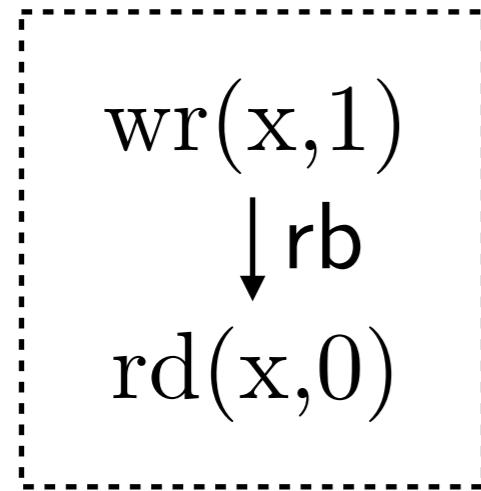


Session Guarantees (Anomalies)

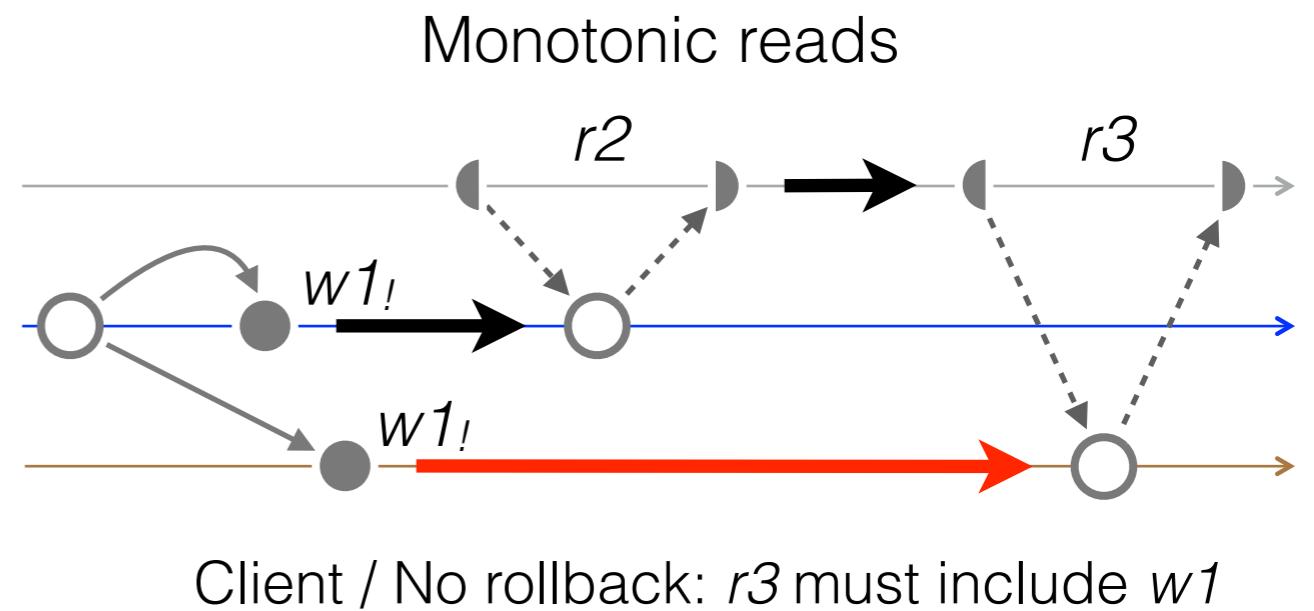
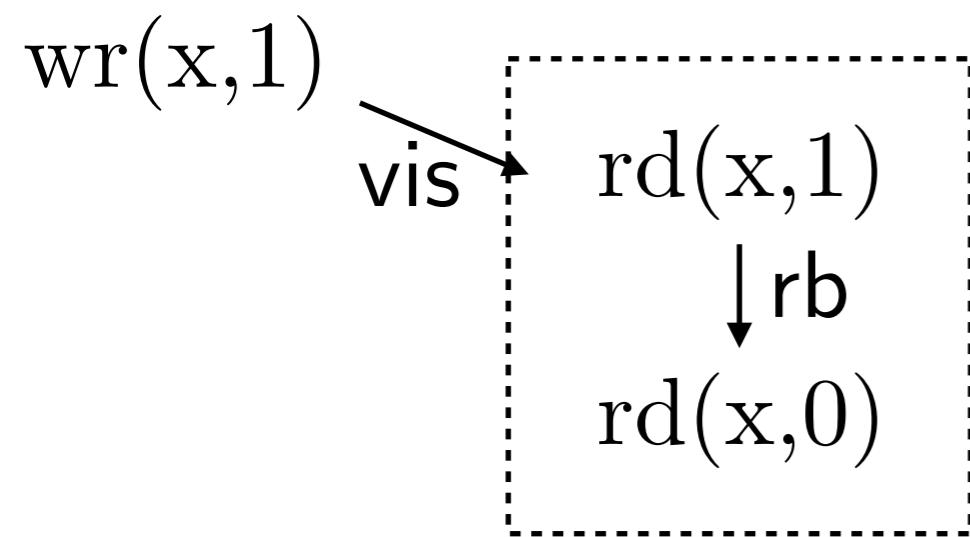
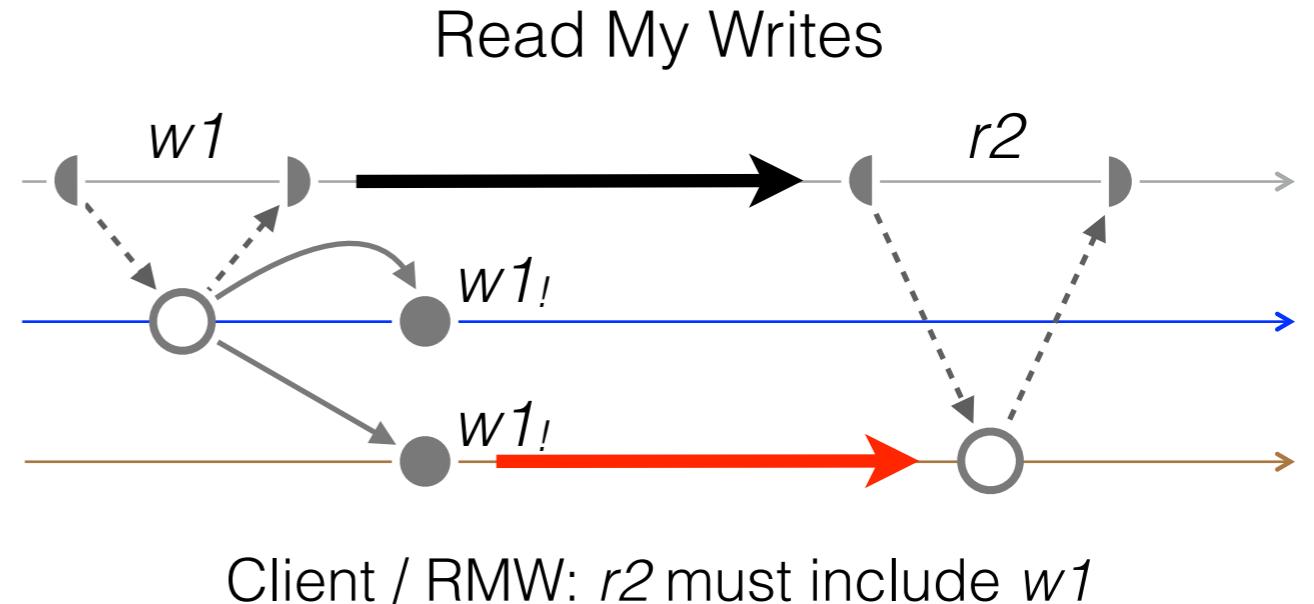
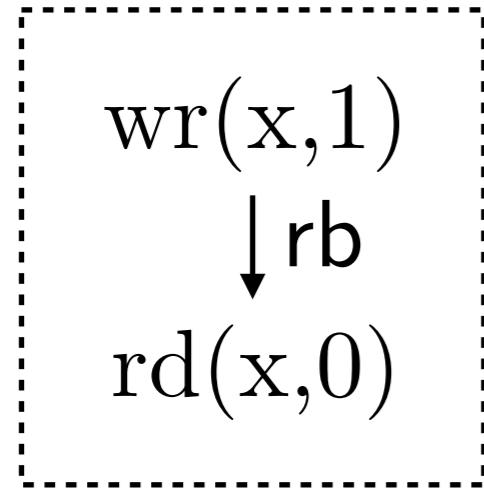
wr(x,1)
↓ rb
rd(x,0)



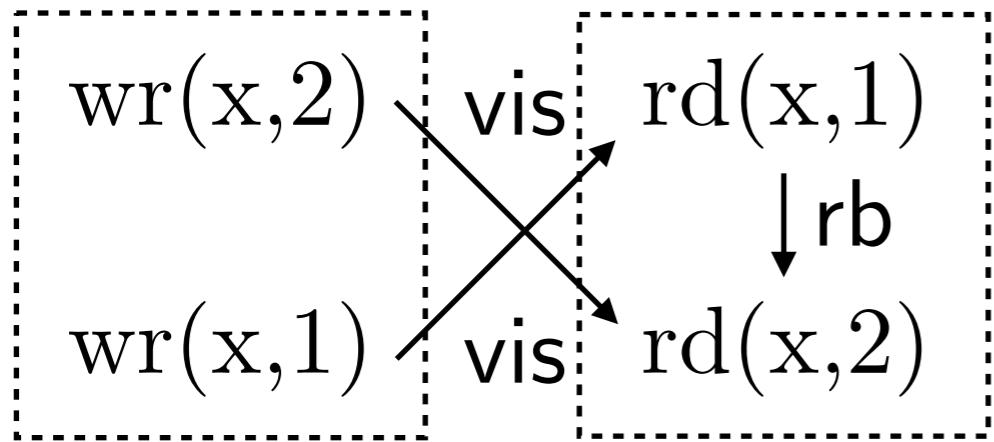
Session Guarantees (Anomalies)



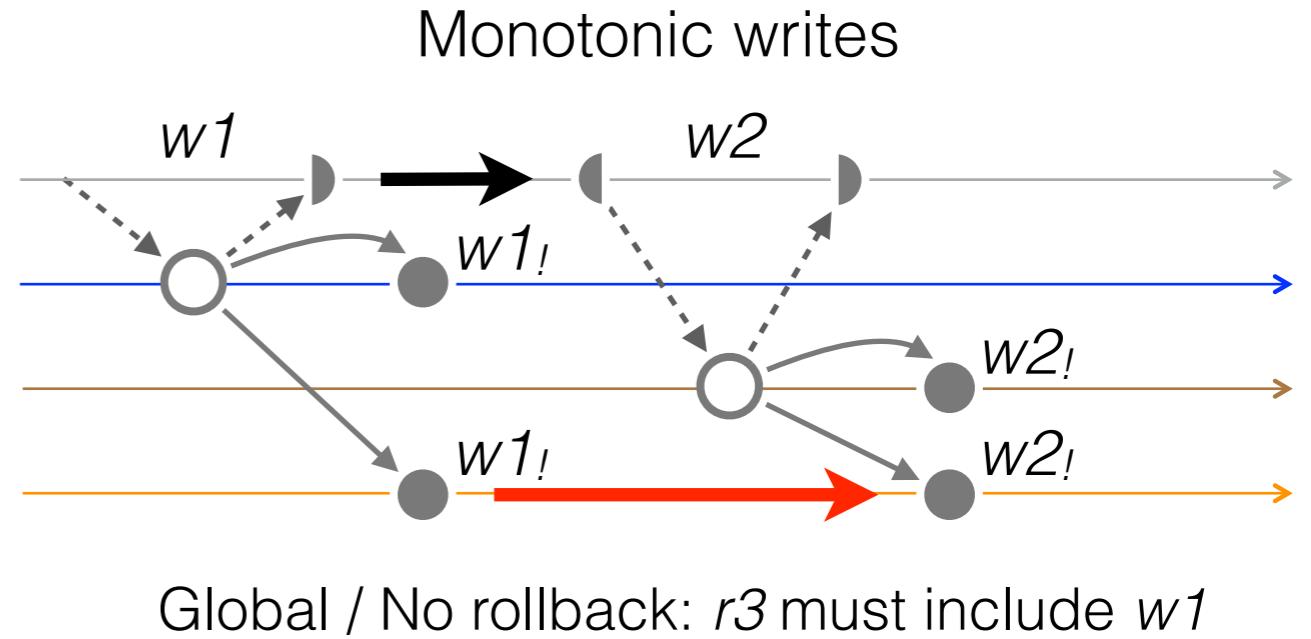
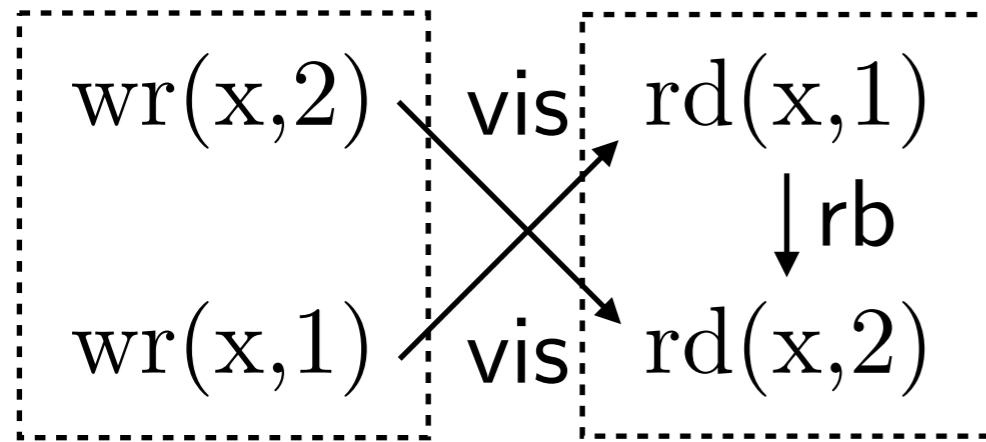
Session Guarantees (Anomalies)



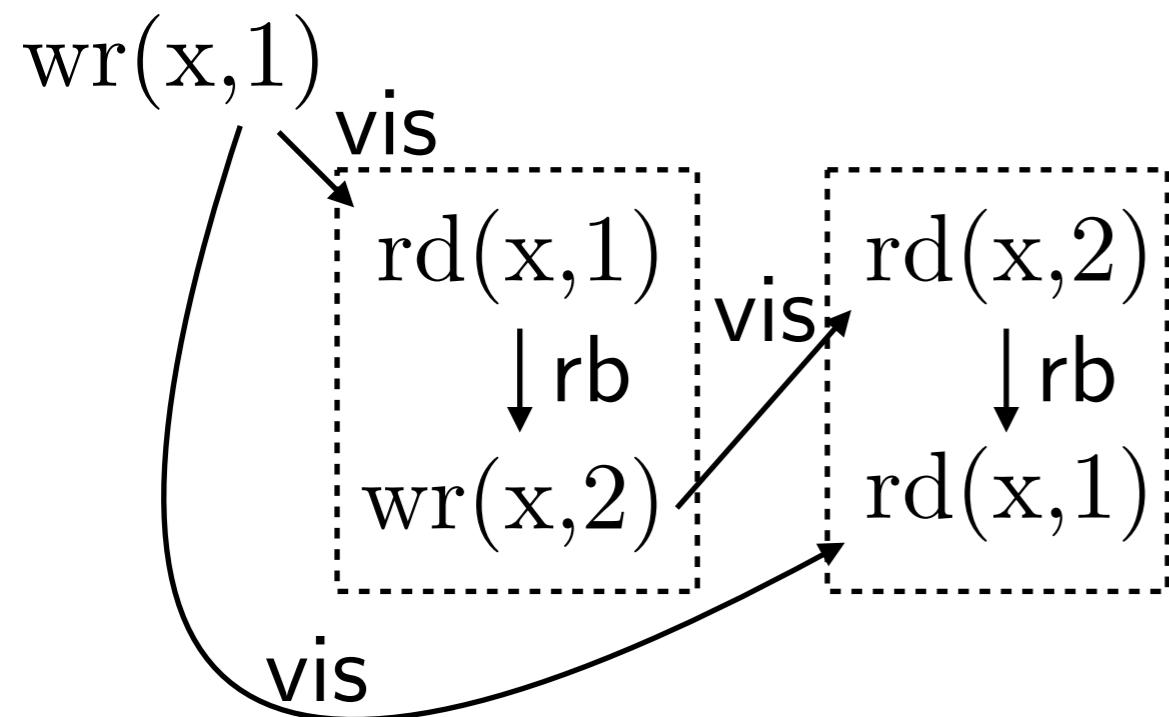
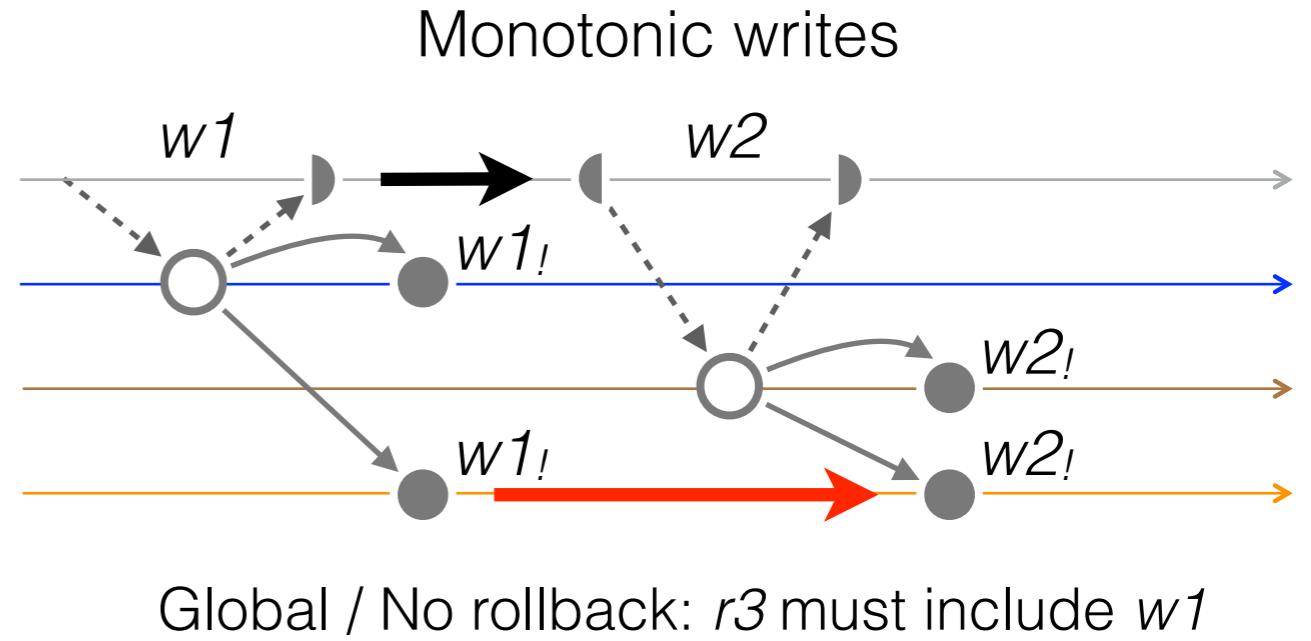
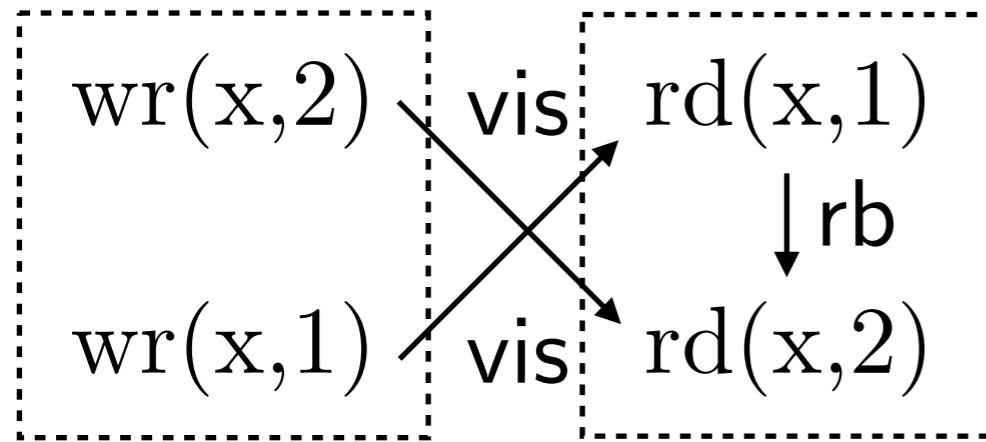
Session Guarantees (Anomalies)



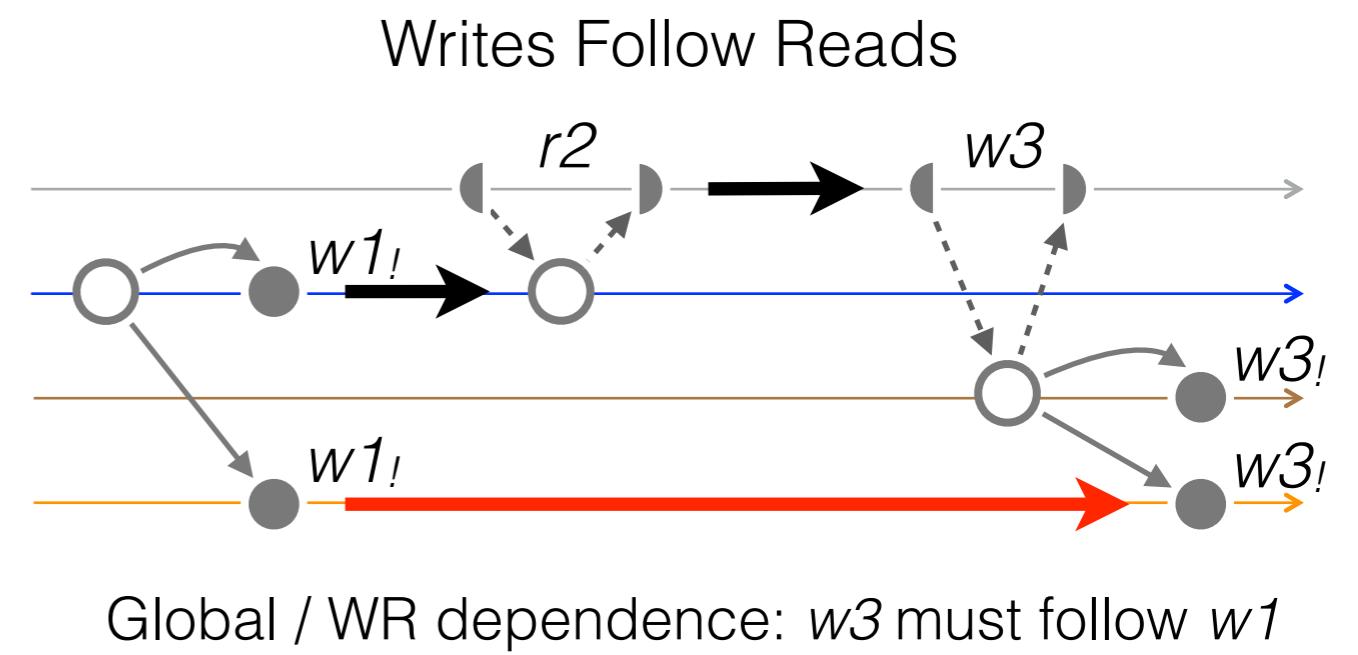
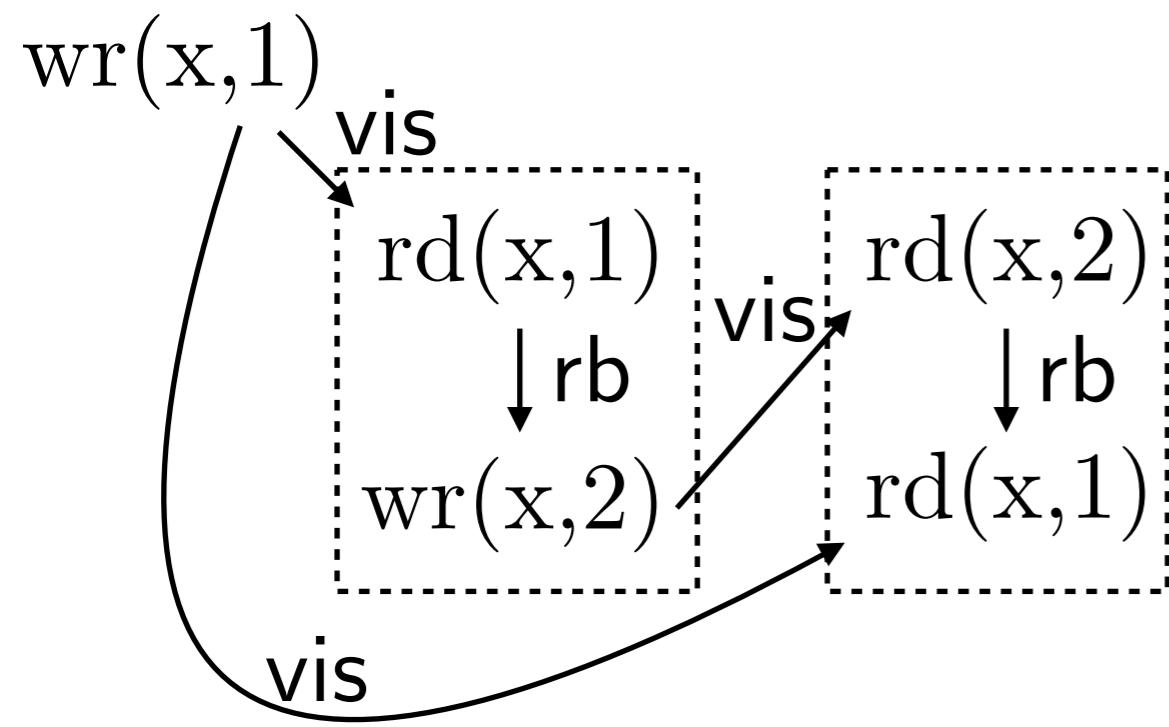
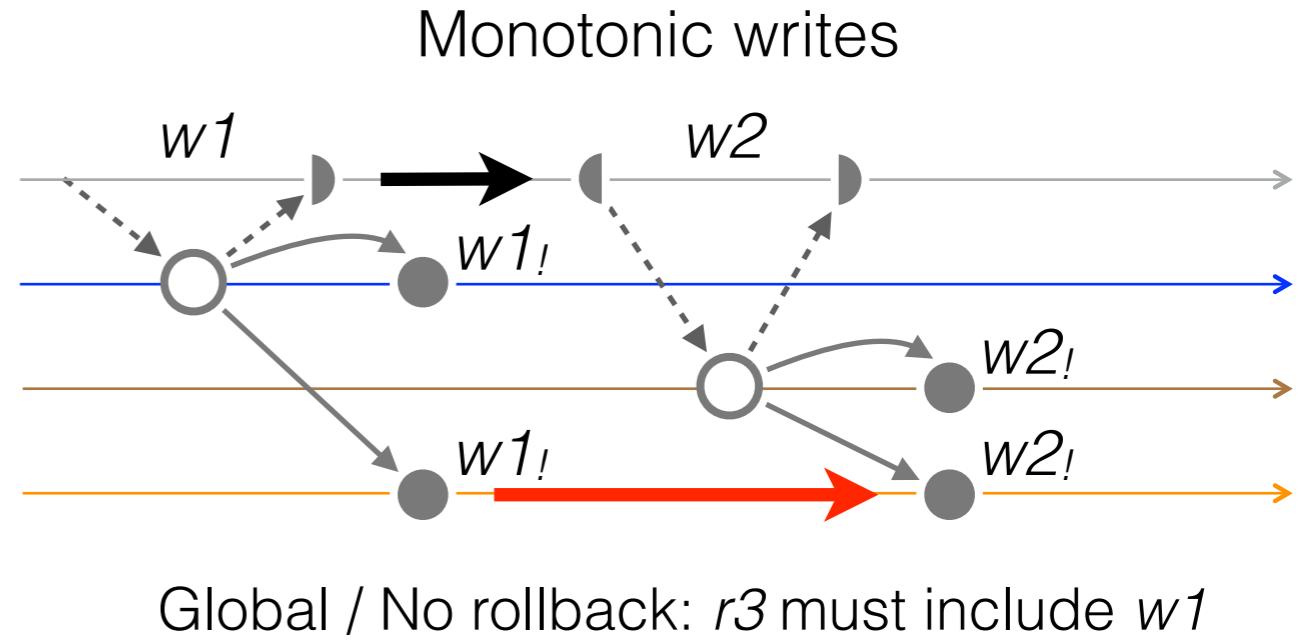
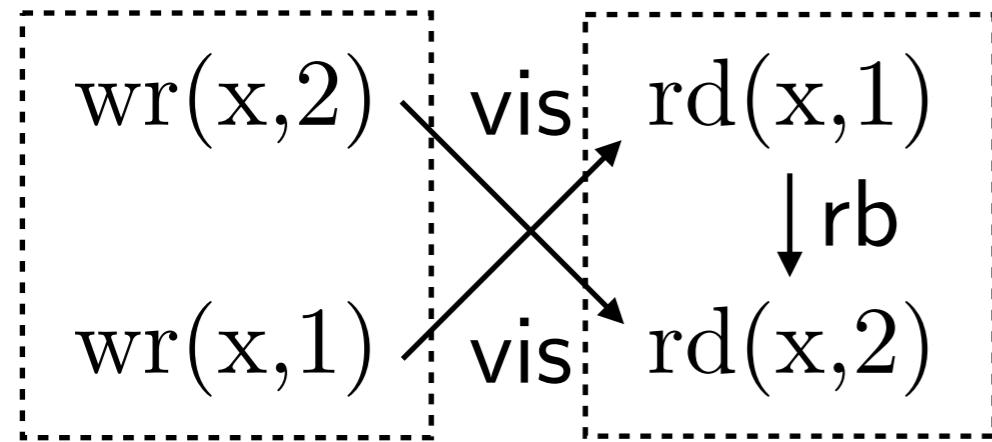
Session Guarantees (Anomalies)



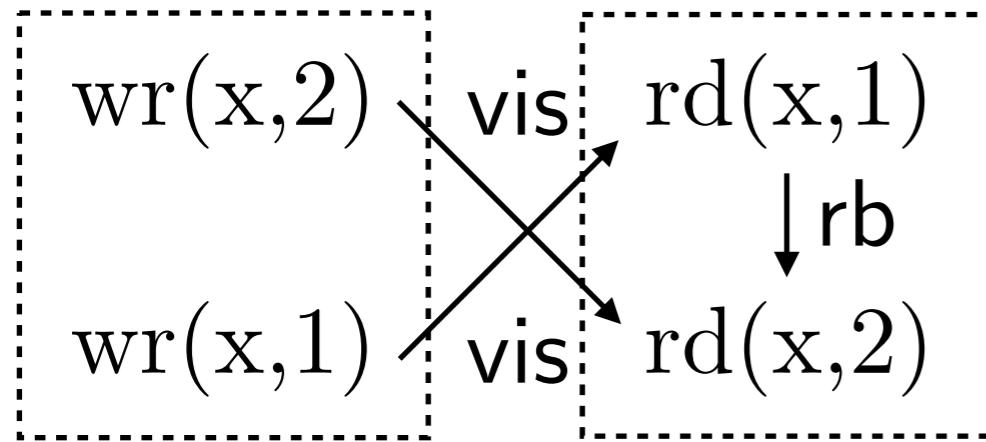
Session Guarantees (Anomalies)



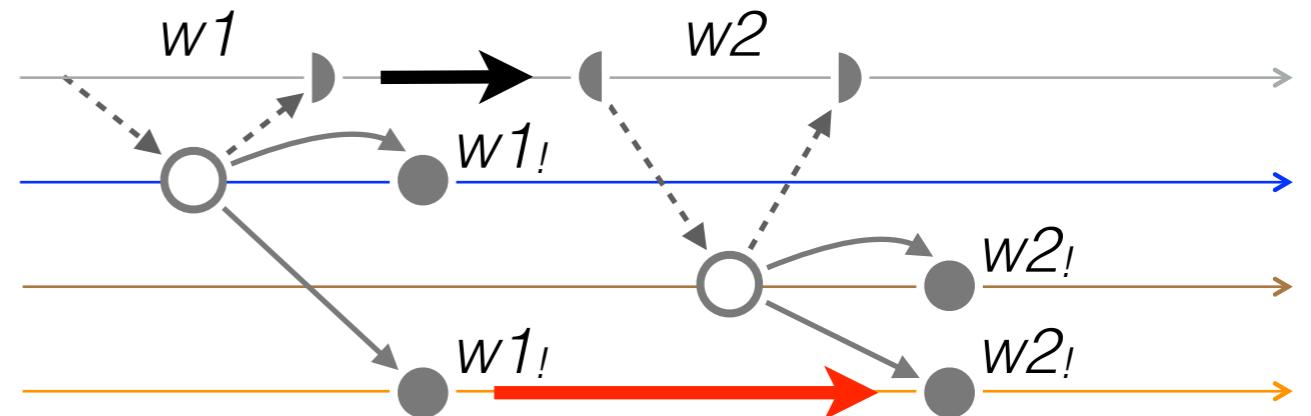
Session Guarantees (Anomalies)



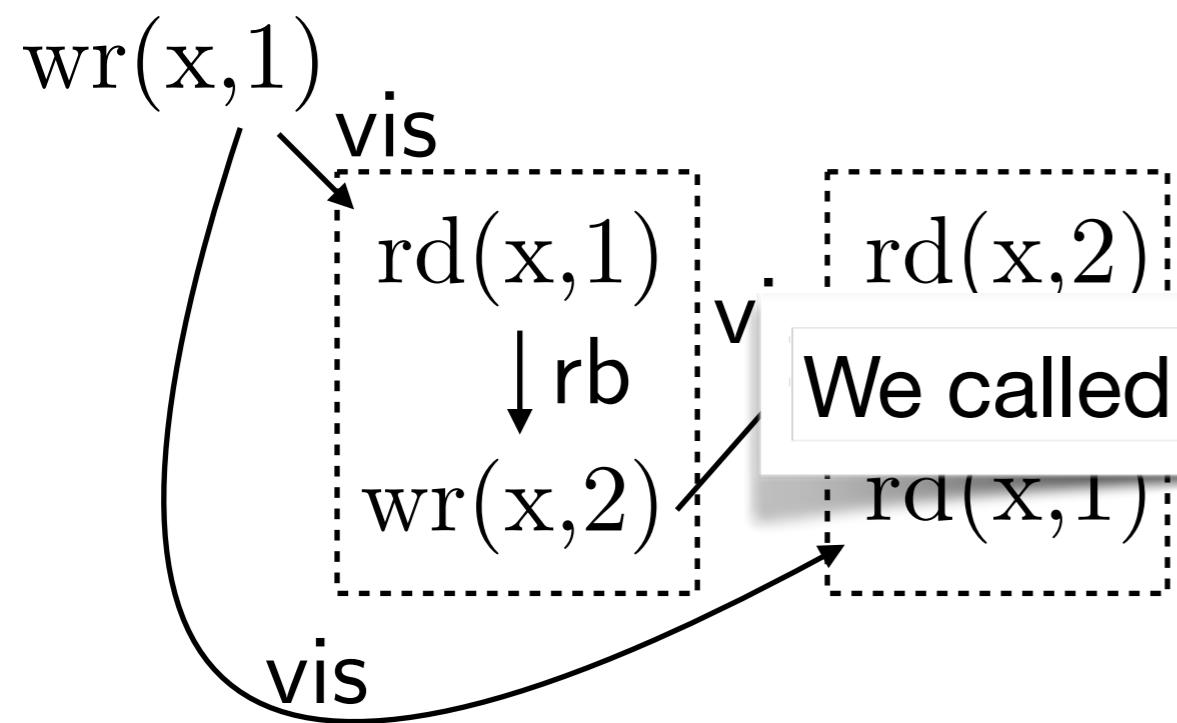
Session Guarantees (Anomalies)



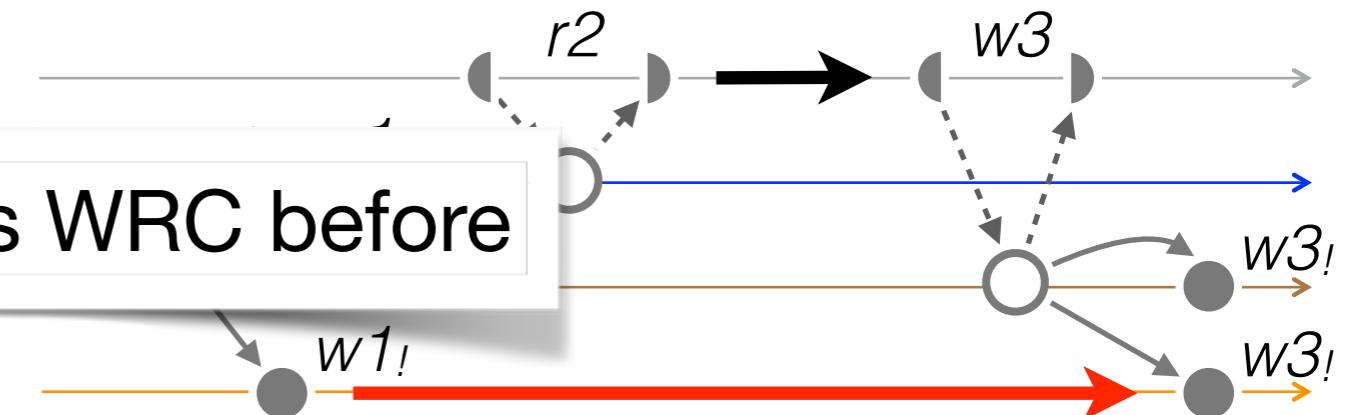
Monotonic writes



Global / No rollback: $r3$ must include $w1$

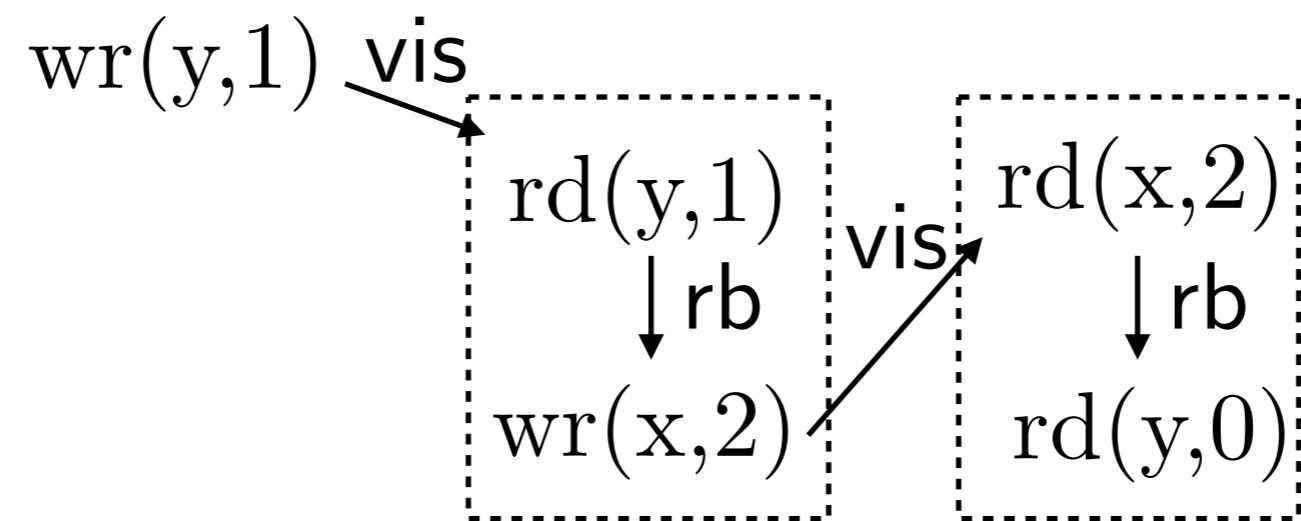


Writes Follow Reads

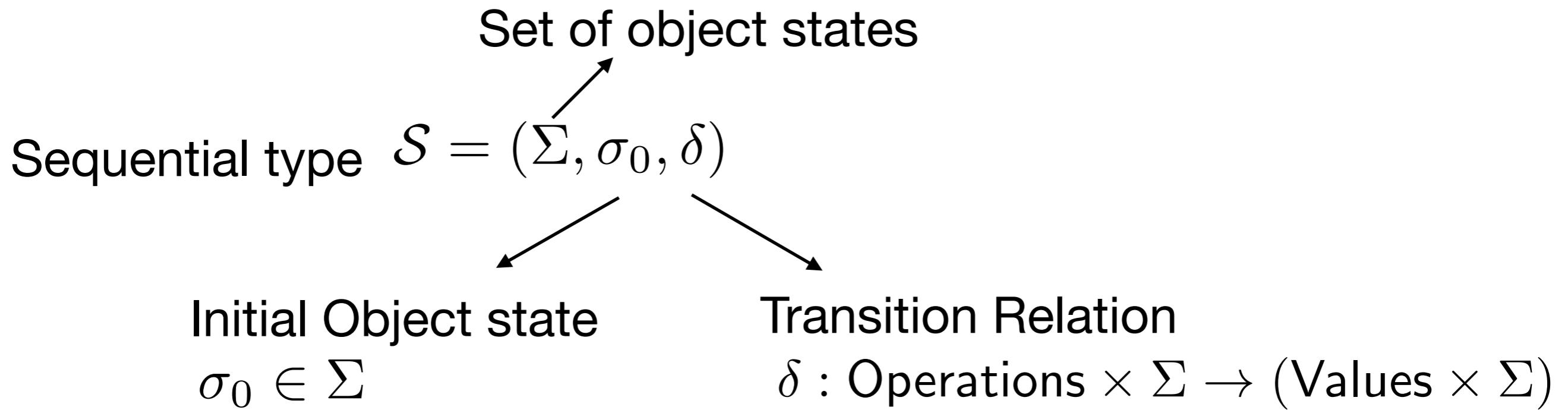


Global / WR dependence: $w3$ must follow $w1$

Causality



Specifying Objects



A register $\mathcal{R}eg = (\mathbb{N}, 0, \delta_r)$

$$\delta_r(n, \text{rd}) = (n, n)$$
$$\delta_r(n, \text{wr}(m)) = (\perp, m)$$

From Burckhardt's

State (and initial state)	Oper.	Returned value	Updated state	Condition
------------------------------	-------	-------------------	------------------	-----------

Counter \mathcal{S}_{ctr}

$n \in \mathbb{N}_0$ (initially 0)	rd	n	same	
	inc	ok	$n + 1$	

Register \mathcal{S}_{reg}

$v \in \text{Values}$ (initially <i>undef</i>)	rd	v	same	
	wr(v')	ok	v'	

Key-Value Store \mathcal{S}_{kvs}

$f : \text{Values} \rightarrow_{\text{fin}}$ Values (initially \emptyset)	rd(k)	<i>undef</i>	same	if $f(k) = \perp$
		$f(k)$	same	if $f(k) \neq \perp$
	wr(k, v)	ok	$f[k \mapsto v]$	

Set \mathcal{S}_{set}

$A \in \mathcal{P}_{\text{fin}}(\text{Values})$ (initially \emptyset)	rd	A	same	
	add(v)	ok	$A \cup \{v\}$	
	rem(v)	ok	$A \setminus \{v\}$	

Queue \mathcal{S}_{queue}

$w \in \text{Values}^*$ (initially ϵ)	enq(v)	ok	$w \cdot v$	
	deq	v	w'	if $w = v \cdot w'$
			∇	if $w = \epsilon$

What is the state in a DS?

- Much like in memory models, there is no unique state σ at every point
- Instead, we have to define the data type based on what is visible at the replica where the operation happens
- We need to change our sequential specifications

Replicated DT Specifications

Instead of a state we use a context for operations

$$C = (E, \text{op}, \text{vis}, \text{ar}) \quad \mathcal{C} \text{ Type of all contexts}$$

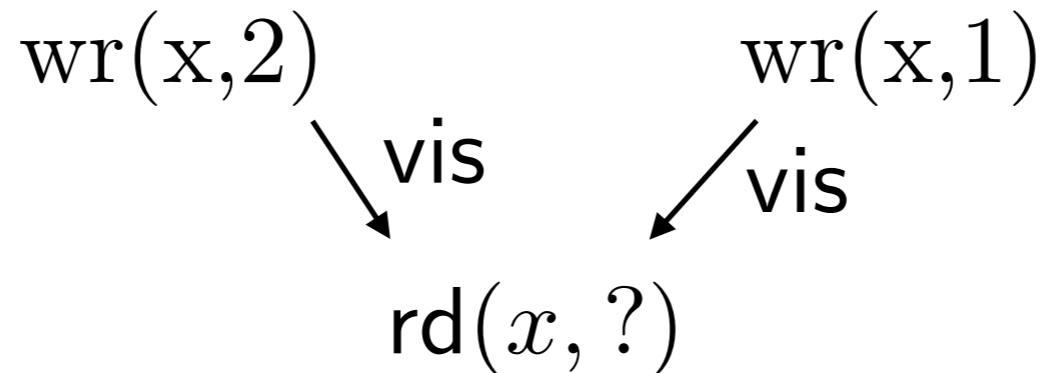
We specify an object based on contexts

$$\mathcal{F} : \text{Operations} \times \mathcal{C} \rightarrow \text{Values}$$

Counter

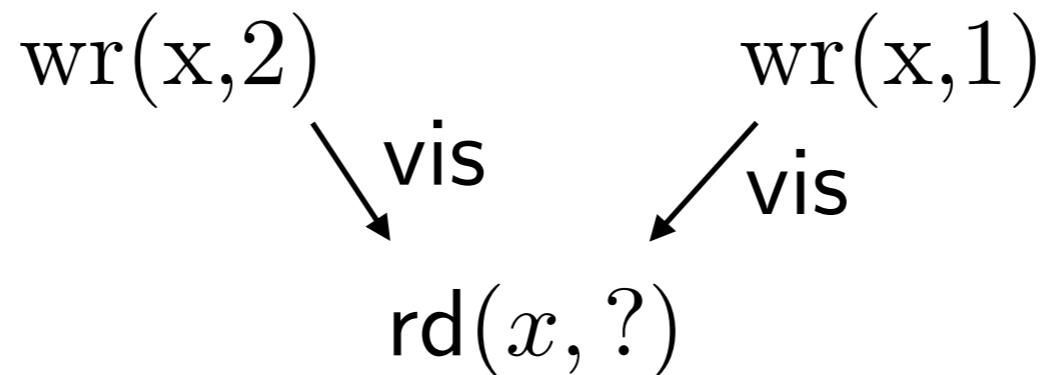
$$\mathcal{F}_{ctr}(\text{rd}, (E, \text{op}, \text{vis}, \text{ar})) = |\{e' \in E \mid \text{op}(e') = \text{inc}\}|$$

What about conflicts



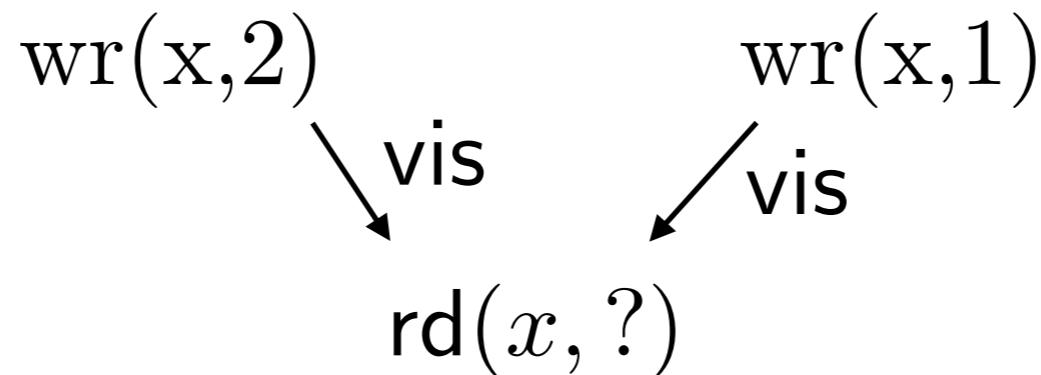
- This non-determinism is problematic

What about conflicts



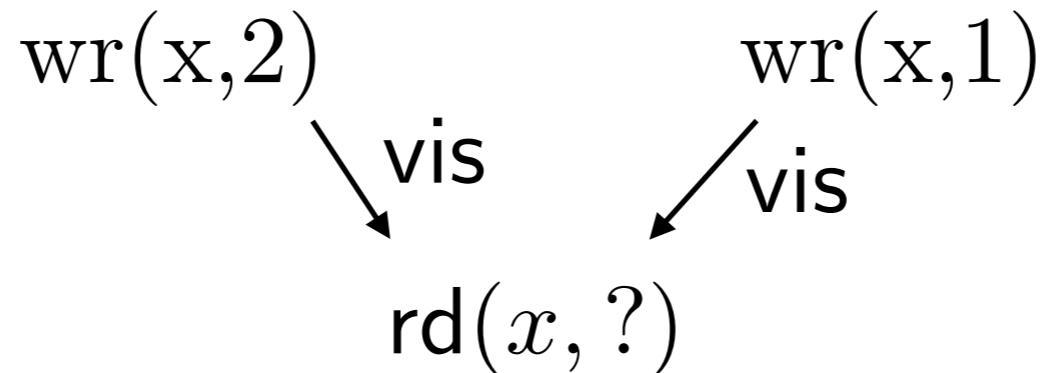
- This non-determinism is problematic
- Could lead to divergent replicas

What about conflicts



- This non-determinism is problematic
- Could lead to divergent replicas
- Yet, synchronization at this scale is too expensive

What about conflicts



- This non-determinism is problematic
- Could lead to divergent replicas
- Yet, synchronization at this scale is too expensive
- Conflict-Free Replicated Data Types to the rescue

What about conflicts



Conflict-free Replicated Data Types *

Marc Shapiro, INRIA & LIP6, Paris, France
Nuno Preguiça, CITI, Universidade Nova de Lisboa, Portugal
Carlos Baquero, Universidade do Minho, Portugal
Marek Zawirski, INRIA & UPMC, Paris, France

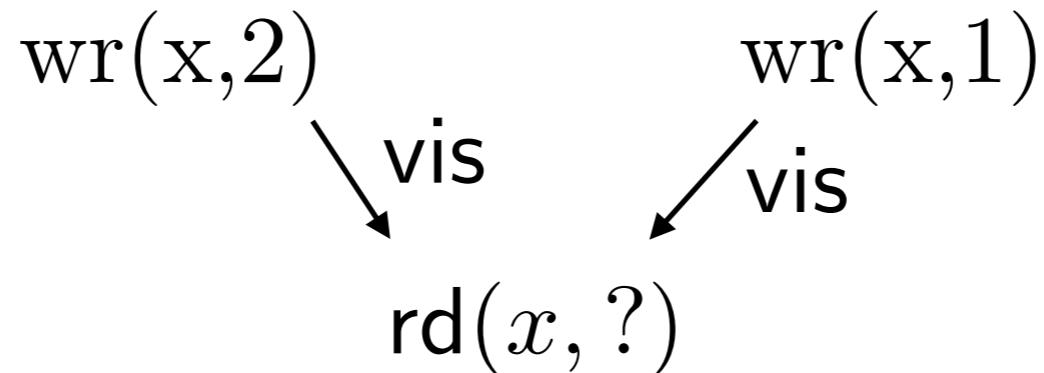
Thème COM — Systèmes communicants
Projet Regal

Rapport de recherche n° 7687 — version 2[†] — version initiale 19 juillet 2011 — version
révisée 25 août 2011 — 18 pages

- This non-determinism
- Could lead to conflicts
- Yet, synchronization
- Conflict-Free Replication

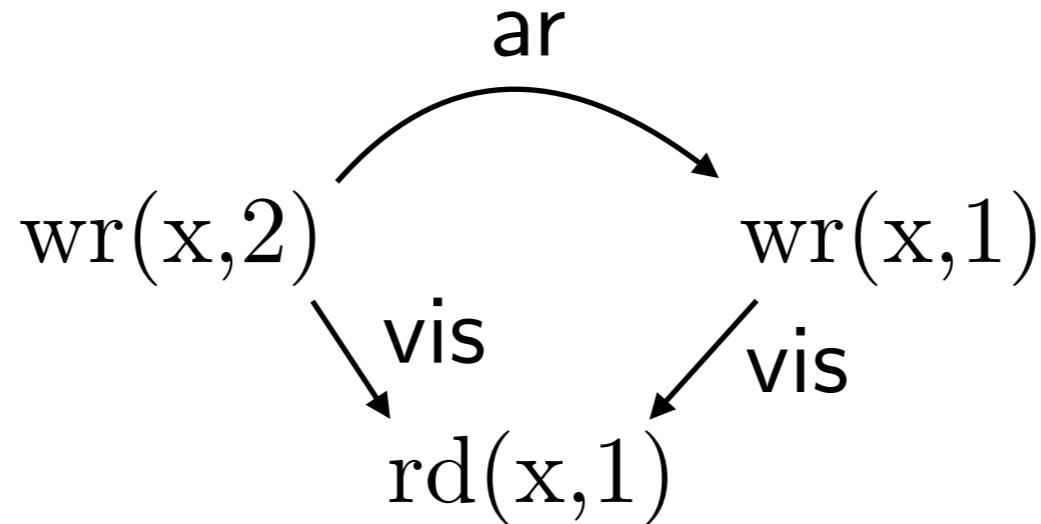
Abstract: Replicating data under Eventual Consistency (EC) allows any replica to accept updates without remote synchronisation. This ensures performance and scalability in large-scale distributed systems (e.g., clouds). However, published EC approaches are ad-hoc and error-prone. Under a formal Strong Eventual Consistency (SEC) model, we study sufficient conditions for convergence. A data type that satisfies these conditions is called a Conflict-free Replicated Data Type (CRDT). Replicas of any CRDT are guaranteed to converge in a self-stabilising manner, despite any number of failures. This paper formalises two popular approaches (state- and operation-based) and their relevant sufficient conditions. We study a complex Graph data type. CRDT types can be used to build interesting theoretical

What about conflicts



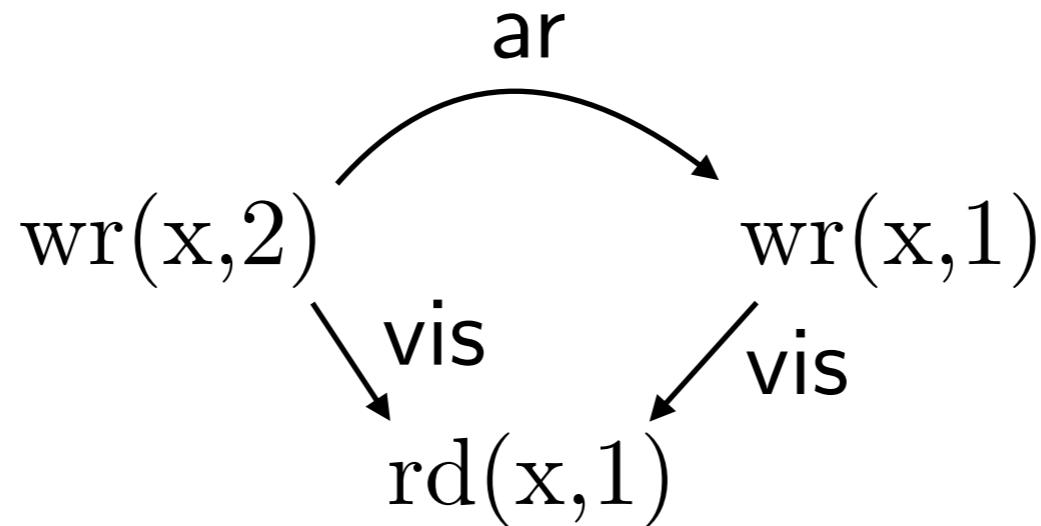
- This non-determinism is problematic
- Could lead to divergent replicas
- Yet, synchronization at this scale is too expensive
- Conflict-Free Replicated Data Types to the rescue
 - aka. Convergent RDTs, Commutative RDTs
 - They enforce a winning strategy between conflicts

What about conflicts



- This non-determinism is problematic
- Could lead to divergent replicas
- Yet, synchronization at this scale is too expensive
- Conflict-Free Replicated Data Types to the rescue
 - aka. Convergent RDTs, Commutative RDTs
 - They enforce a winning strategy between conflicts

What about conflicts



Last Writer Wins Register

$$\mathcal{F}_{\text{reg}}(\text{rd}, (E, \text{op}, \text{vis}, \text{ar})) = \begin{cases} \text{undef} & \text{if } \text{writes}(E) = \emptyset \\ v & \text{if } \text{op}(\max_{\text{ar}} \text{writes}(E)) = \text{wr}(v) \end{cases}$$

RDTs

Multi-Value Register

$$\mathcal{F}_{mvr}(\text{rd}, (E, \text{op}, \text{vis}, \text{ar})) =$$

$$\{v \mid \exists e \in E : \text{op}(e) = \text{wr}(v) \text{ and } \forall e' \in \text{writes}(E) : e \not\overset{\text{vis}}{\rightarrow} e'\}$$

RDTs

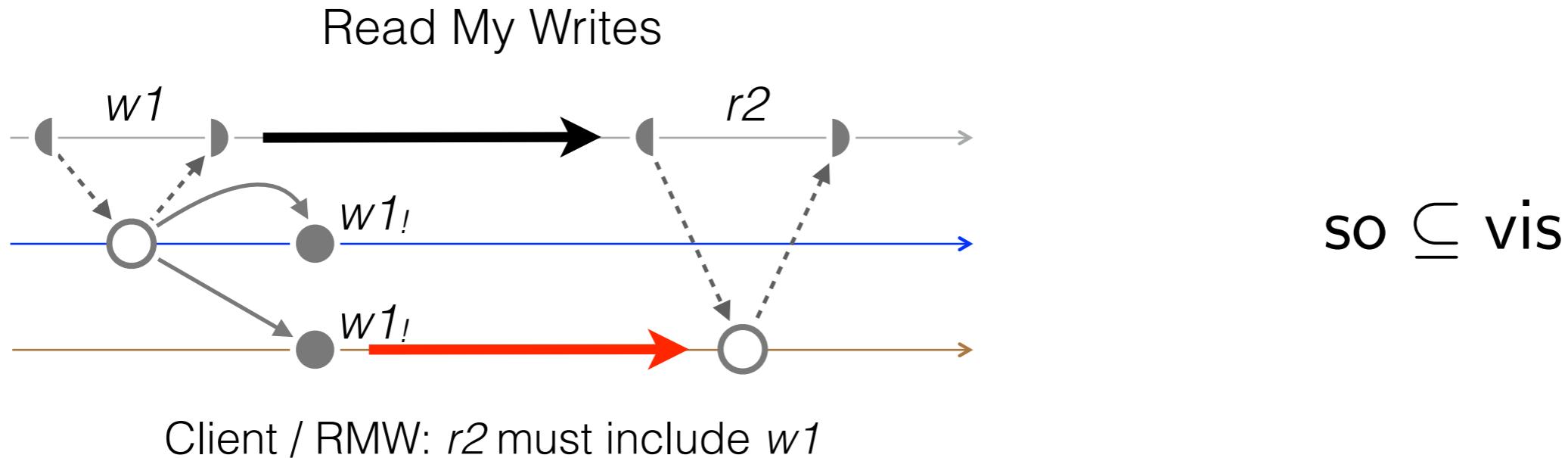
Multi-Value Register

$$\begin{aligned}\mathcal{F}_{mvr}(\text{rd}, (E, \text{op}, \text{vis}, \text{ar})) = \\ \{v \mid \exists e \in E : \text{op}(e) = \text{wr}(v) \text{ and } \forall e' \in \text{writes}(E) : e \xrightarrow[\text{vis}]{} e'\}\end{aligned}$$

Add-wins set

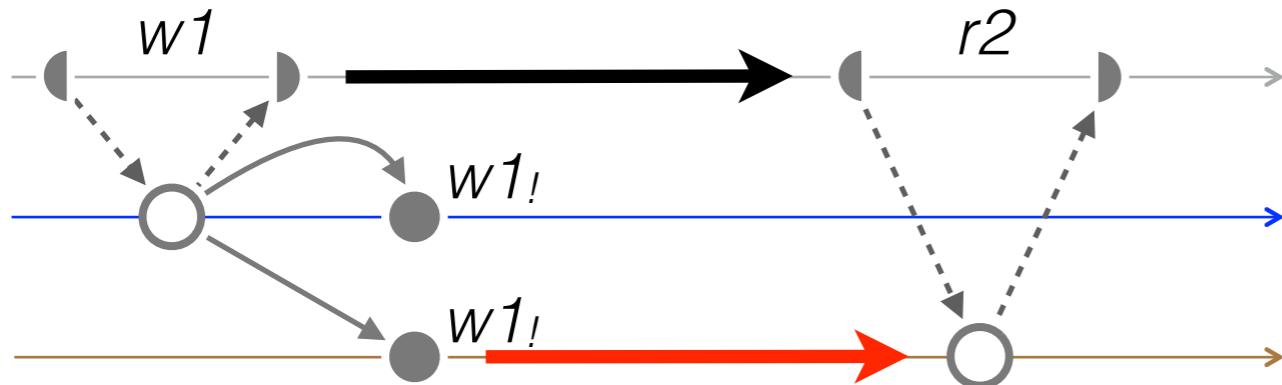
$$\begin{aligned}\mathcal{F}_{awset}(\text{contains}(v), (E, \text{op}, \text{vis}, \text{ar})) = \text{true} &\iff \\ \exists e \in E : \text{op}(e) = \text{add}(v) \wedge \neg(\exists e' \in E : \text{op}(e') = \text{rem}(v) \wedge e \xrightarrow{\text{vis}} e')\end{aligned}$$

Consistency Axioms



Consistency Axioms

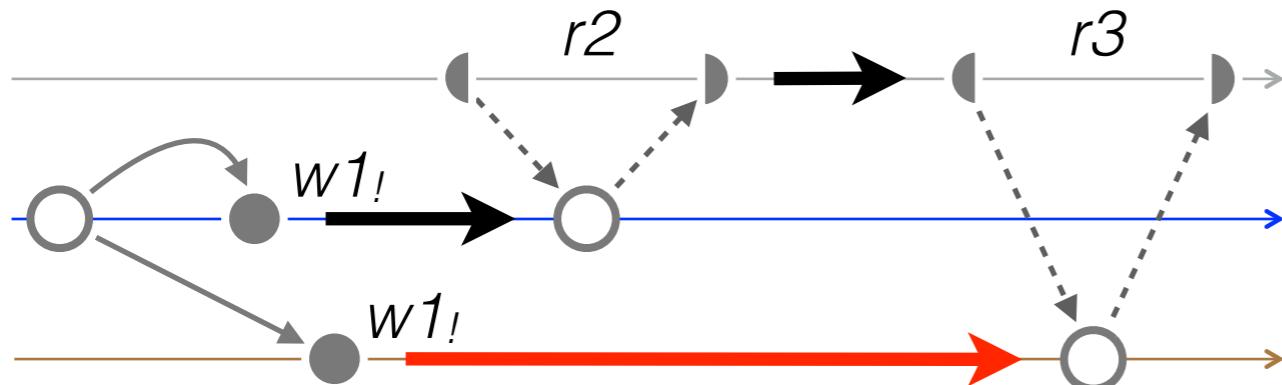
Read My Writes



$\text{so} \subseteq \text{vis}$

Client / RMW: $r2$ must include $w1$

Monotonic reads

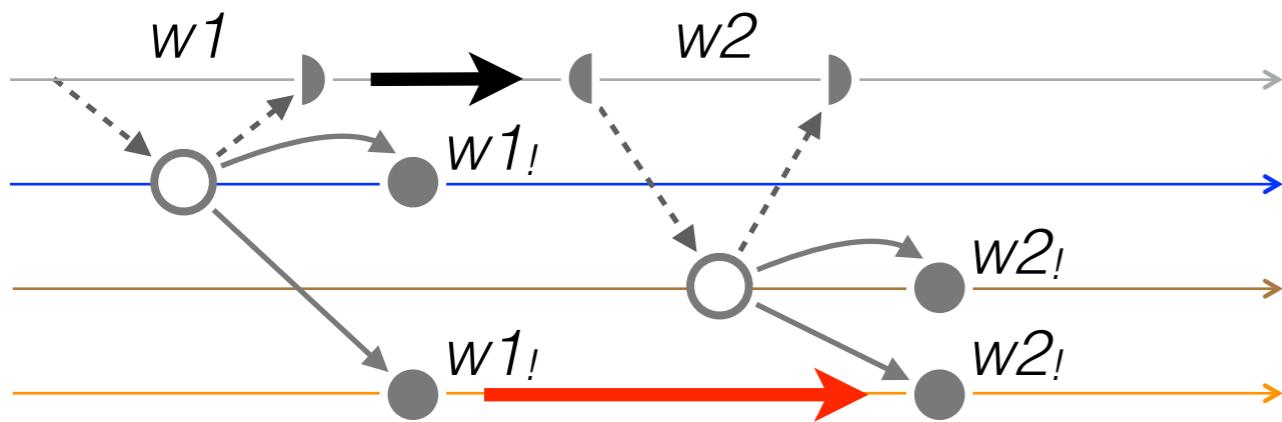


$\text{vis}; \text{so} \subseteq \text{vis}$

Client / No rollback: $r3$ must include $w1$

Consistency Axioms

Monotonic writes

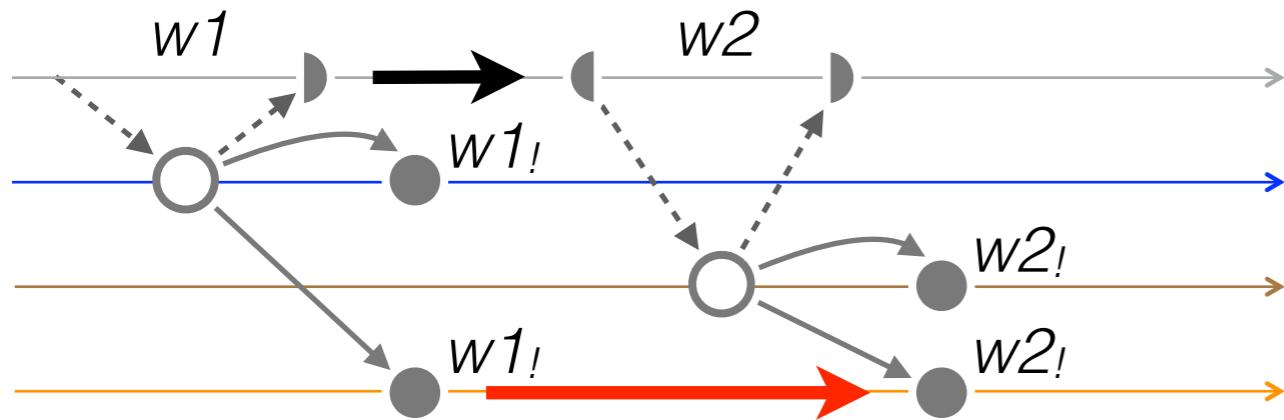


$$ss \cap (wr \times wr) \subseteq ar$$

Global / No rollback: $r3$ must include $w1$

Consistency Axioms

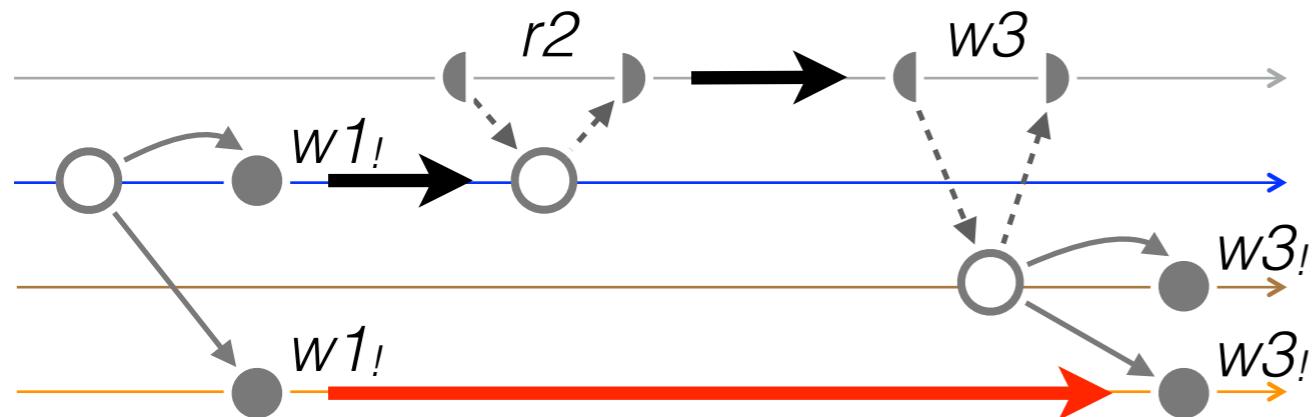
Monotonic writes



$$ss \cap (wr \times wr) \subseteq ar$$

Global / No rollback: $r3$ must include $w1$

Writes Follow Reads



Causal Visibility
 $hb = (vis \cup so)^+$

$$hb \subseteq vis$$

Global / WR dependence: $w3$ must follow $w1$

Consistency Models

$$\text{READMYWRITES} \stackrel{\text{def}}{=} (\text{so} \subseteq \text{vis})$$

$$\text{MONOTONICREADS} \stackrel{\text{def}}{=} (\text{vis} ; \text{so}) \subseteq \text{vis}$$

$$\text{CONSISTENTPREFIX} \stackrel{\text{def}}{=} (\text{ar} ; (\text{vis} \cap \neg \text{ss})) \subseteq \text{vis}$$

$$\text{NO CIRCULAR CAUSALITY} \stackrel{\text{def}}{=} \text{acyclic}(\text{hb})$$

$$\text{CAUSALVISIBILITY} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{vis})$$

$$\text{CAUSALARBITRATION} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{ar})$$

$$\text{CAUSALITY} \stackrel{\text{def}}{=} \text{CAUSALVISIBILITY} \wedge \text{CAUSALARBITRATION}$$

$$\text{SINGLEORDER} \stackrel{\text{def}}{=} \text{vis} = \text{ar}$$

$$\text{REALTIME} \stackrel{\text{def}}{=} \text{rb} \subseteq \text{ar}$$

Consistency Models

$$\text{LINEARIZABILITY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{REALTIME} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{READMYWRITES} \stackrel{\text{def}}{=} (\text{so} \subseteq \text{vis})$$

$$\text{MONOTONICREADS} \stackrel{\text{def}}{=} (\text{vis} ; \text{so}) \subseteq \text{vis}$$

$$\text{CONSISTENTPREFIX} \stackrel{\text{def}}{=} (\text{ar} ; (\text{vis} \cap \neg \text{ss})) \subseteq \text{vis}$$

$$\text{NO CIRCULAR CAUSALITY} \stackrel{\text{def}}{=} \text{acyclic}(\text{hb})$$

$$\text{CAUSALVISIBILITY} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{vis})$$

$$\text{CAUSALARBITRATION} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{ar})$$

$$\text{CAUSALITY} \stackrel{\text{def}}{=} \text{CAUSALVISIBILITY} \wedge \text{CAUSALARBITRATION}$$

$$\text{SINGLEORDER} \stackrel{\text{def}}{=} \text{vis} = \text{ar}$$

$$\text{REALTIME} \stackrel{\text{def}}{=} \text{rb} \subseteq \text{ar}$$

Consistency Models

$$\text{LINEARIZABILITY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{REALTIME} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{SEQUENTIAL CONSISTENCY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{READMYWRITES} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{READMYWRITES} \stackrel{\text{def}}{=} (\text{so} \subseteq \text{vis})$$

$$\text{MONOTONICREADS} \stackrel{\text{def}}{=} (\text{vis} ; \text{so}) \subseteq \text{vis}$$

$$\text{CONSISTENTPREFIX} \stackrel{\text{def}}{=} (\text{ar} ; (\text{vis} \cap \neg \text{ss})) \subseteq \text{vis}$$

$$\text{NO CIRCULAR CAUSALITY} \stackrel{\text{def}}{=} \text{acyclic}(\text{hb})$$

$$\text{CAUSALVISIBILITY} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{vis})$$

$$\text{CAUSALARBITRATION} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{ar})$$

$$\text{CAUSALITY} \stackrel{\text{def}}{=} \text{CAUSALVISIBILITY} \wedge \text{CAUSALARBITRATION}$$

$$\text{SINGLEORDER} \stackrel{\text{def}}{=} \text{vis} = \text{ar}$$

$$\text{REALTIME} \stackrel{\text{def}}{=} \text{rb} \subseteq \text{ar}$$

Consistency Models

$$\text{LINEARIZABILITY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{REALTIME} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{SEQUENTIAL CONSISTENCY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{READMYWRITES} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{CAUSAL CONSISTENCY}(\mathcal{F}) = \text{CAUSALITY} \wedge \text{RVAL}(\mathcal{F})$$

$$\text{READMYWRITES} \stackrel{\text{def}}{=} (\text{so} \subseteq \text{vis})$$

$$\text{MONOTONICREADS} \stackrel{\text{def}}{=} (\text{vis} ; \text{so}) \subseteq \text{vis}$$

$$\text{CONSISTENTPREFIX} \stackrel{\text{def}}{=} (\text{ar} ; (\text{vis} \cap \neg \text{ss})) \subseteq \text{vis}$$

$$\text{NO CIRCULAR CAUSALITY} \stackrel{\text{def}}{=} \text{acyclic}(\text{hb})$$

$$\text{CAUSALVISIBILITY} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{vis})$$

$$\text{CAUSALARBITRATION} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{ar})$$

$$\text{CAUSALITY} \stackrel{\text{def}}{=} \text{CAUSALVISIBILITY} \wedge \text{CAUSALARBITRATION}$$

$$\text{SINGLEORDER} \stackrel{\text{def}}{=} \text{vis} = \text{ar}$$

$$\text{REALTIME} \stackrel{\text{def}}{=} \text{rb} \subseteq \text{ar}$$

Consistency Models

$\text{LINEARIZABILITY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{REALTIME} \wedge \text{RVAL}(\mathcal{F})$

$\text{SEQUENTIAL CONSISTENCY}(\mathcal{F}) = \text{SINGLEORDER} \wedge \text{READMYWRITES} \wedge \text{RVAL}(\mathcal{F})$

$\text{CAUSAL CONSISTENCY}(\mathcal{F}) = \text{CAUSALITY} \wedge \text{RVAL}(\mathcal{F})$

$\text{EVENTUAL CONSISTENCY}(\mathcal{F}) = \text{NoCIRCULARCAUSALITY} \wedge \text{RVAL}(\mathcal{F})$

$$\text{READMYWRITES} \stackrel{\text{def}}{=} (\text{so} \subseteq \text{vis})$$

$$\text{MONOTONICREADS} \stackrel{\text{def}}{=} (\text{vis} ; \text{so}) \subseteq \text{vis}$$

$$\text{CONSISTENTPREFIX} \stackrel{\text{def}}{=} (\text{ar} ; (\text{vis} \cap \neg \text{ss})) \subseteq \text{vis}$$

$$\text{NoCIRCULARCAUSALITY} \stackrel{\text{def}}{=} \text{acyclic}(\text{hb})$$

$$\text{CAUSALVISIBILITY} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{vis})$$

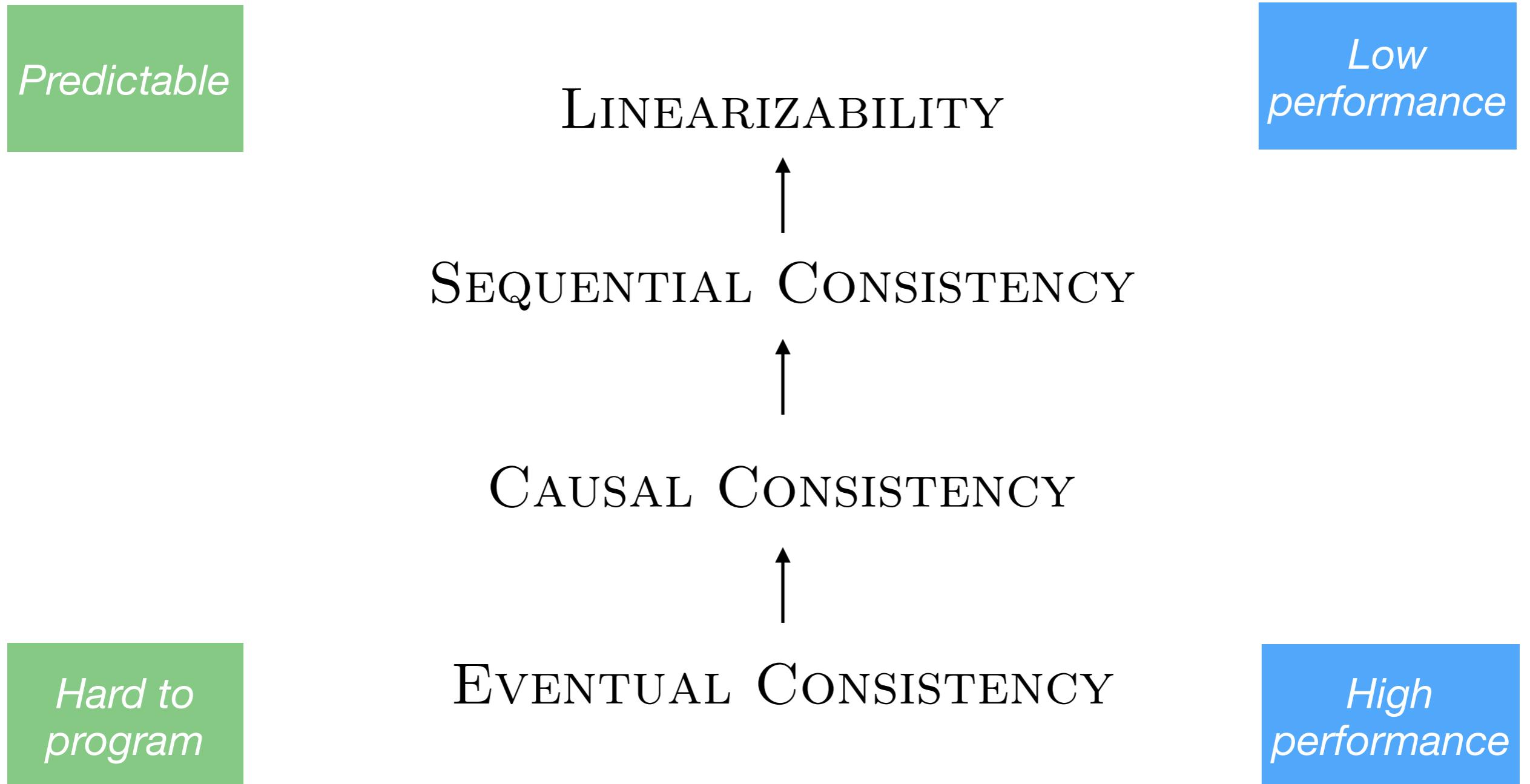
$$\text{CAUSALARBITRATION} \stackrel{\text{def}}{=} (\text{hb} \subseteq \text{ar})$$

$$\text{CAUSALITY} \stackrel{\text{def}}{=} \text{CAUSALVISIBILITY} \wedge \text{CAUSALARBITRATION}$$

$$\text{SINGLEORDER} \stackrel{\text{def}}{=} \text{vis} = \text{ar}$$

$$\text{REALTIME} \stackrel{\text{def}}{=} \text{rb} \subseteq \text{ar}$$

Strong vs. weak?



Replicated Data Bases

Replicated Data Bases

A Framework for Transactional Consistency Models with Atomic Visibility

Andrea Cerone, Giovanni Bernardi, and Alexey Gotsman

IMDEA Software Institute, Madrid, Spain

Abstract

Modern distributed systems often rely on databases that achieve scalability by providing only weak guarantees about the consistency of distributed transaction processing. The semantics of programs interacting with such a database depends on its consistency model, defining these guarantees. Unfortunately, consistency models are usually stated informally or using disparate formalisms, often tied to the database internals. To deal with this problem, we propose a framework for specifying a variety of consistency models for transactions uniformly and declaratively. Our specifications are given in the style of weak memory models, using structures of events and relations on them. The specifications are particularly concise because they exploit the property of atomic visibility guaranteed by many consistency models: either all or none of the updates by a transaction can be visible to another one. This allows the specifications to abstract from individual events inside transactions. We illustrate the use of our framework by specifying several existing consistency models. To validate our specifications, we prove that they are equivalent to alternative operational ones, given as algorithms closer to actual implementations. Our work provides a rigorous foundation for developing the metatheory of the novel form of concurrency arising in weakly consistent large-scale databases.

1998 ACM Subject Classification C.2.4 Distributed Systems

Keywords and phrases Replication, Consistency models, Transactions

Digital Object Identifier 10.4230/LIPIcs.xxx.yyy.p

1 Introduction

To achieve availability and scalability, modern distributed systems often rely on *replicated databases*, which maintain multiple *replicas* of shared data. The database clients can execute transactions on the data at any of the replicas, which communicate changes to each other using message passing. For example, large-scale Internet services use data replicas in geographically distinct locations, and applications for mobile devices keep replicas locally as well as in the cloud to support offline use. Ideally, we want the concurrent and distributed processing in a replicated database to be transparent, as formalised by the classical notion of *atomicity*: the database behaves as if it executed transactions serially on a non-distributed database. In ideal requires extensive coordination between the database instances and the network.

Replicated Data Bases

- Operations are transactions
- Each transaction issues a number of reads and writes
- Writes are ordered by program order (po)
- Visibility and Arbitration relate “transactions” instead of individual reads and writes
- We consider only the case of a key-value store data base

Replicated Data Bases

Replicated Data Bases

Transaction $T = (E, \text{po})$

Replicated Data Bases

Transaction $T = (E, \text{po})$

History $H = \{T_0, T_1, \dots, T_n\}$

Replicated Data Bases

Transaction $T = (E, \text{po})$

History $H = \{T_0, T_1, \dots, T_n\}$

Abstract Execution $A = (H, \text{vis}, \text{ar})$

Replicated Data Bases

Transaction $T = (E, \text{po})$

History $H = \{T_0, T_1, \dots, T_n\}$

Abstract Execution $A = (H, \text{vis}, \text{ar})$

These are relation on transactions now

$$T_1 \xrightarrow{\text{vis}} T_2$$

Replicated Data Bases

Consistency Axioms

Replicated Data Bases

Consistency Axioms

Read Consistency

- Every read in transaction T sees a visible write to T that is the maximum visible write according to arbitration

$$rd(x, v) \in E(T) \Rightarrow \exists T', \quad T' \in \text{vis}_H^{-1}(T) \wedge \\ wr(x, 1) \in E(T') \wedge \\ \max_{\text{ar}}(\text{vis}_H^{-1}(T) \cap Writes_H(x)) = T'$$

Replicated Data Bases

Consistency Axioms

Read Consistency

- Every read in transaction T sees a visible write to T that is the maximum visible write according to arbitration

$$rd(x, v) \in E(T) \Rightarrow \exists T', \quad T' \in \text{vis}_H^{-1}(T) \wedge \\ wr(x, 1) \in E(T') \wedge \\ \max_{\text{ar}}(\text{vis}_H^{-1}(T) \cap Writes_H(x)) = T'$$

Transitive Visibility $\text{vis}^+ \subseteq \text{vis}$

Replicated Data Bases

Consistency Axioms

Read Consistency

- Every read in transaction T sees a visible write to T that is the maximum visible write according to arbitration

$$rd(x, v) \in E(T) \Rightarrow \exists T', \quad T' \in \text{vis}_H^{-1}(T) \wedge \\ wr(x, 1) \in E(T') \wedge \\ \max_{\text{ar}}(\text{vis}_H^{-1}(T) \cap \text{Writes}_H(x)) = T'$$

Transitive Visibility $\text{vis}^+ \subseteq \text{vis}$

Prefix Consistent $\text{ar}; \text{vis} \subseteq \text{vis}$

Replicated Data Bases

Consistency Axioms

Read Consistency

- Every read in transaction T sees a visible write to T that is the maximum visible write according to arbitration

$$rd(x, v) \in E(T) \Rightarrow \exists T', \quad T' \in \text{vis}_H^{-1}(T) \wedge \\ wr(x, 1) \in E(T') \wedge \\ \max_{\text{ar}}(\text{vis}_H^{-1}(T) \cap \text{Writes}_H(x)) = T'$$

Transitive Visibility $\text{vis}^+ \subseteq \text{vis}$

Prefix Consistent $\text{ar}; \text{vis} \subseteq \text{vis}$

Total Visibility $\forall T_1, T_2. \quad T_1 \xrightarrow{\text{vis}} T_2 \vee T_2 \xrightarrow{\text{vis}} T_1$

Replicated Data Bases

Consistency Axioms

Read Consistency

- Every read in transaction T sees a visible write to T that is the maximum visible write according to arbitration

$$rd(x, v) \in E(T) \Rightarrow \exists T', \quad T' \in \text{vis}_H^{-1}(T) \wedge \\ wr(x, 1) \in E(T') \wedge \\ \max_{\text{ar}}(\text{vis}_H^{-1}(T) \cap \text{Writes}_H(x)) = T'$$

Transitive Visibility $\text{vis}^+ \subseteq \text{vis}$

Prefix Consistent $\text{ar}; \text{vis} \subseteq \text{vis}$

Total Visibility $\forall T_1, T_2. \quad T_1 \xrightarrow{\text{vis}} T_2 \vee T_2 \xrightarrow{\text{vis}} T_1$

No Conflict $\{T_1, T_2\} \subseteq \text{Writes}_H(x) \Rightarrow T_1 \xrightarrow{\text{vis}} T_2 \vee T_2 \xrightarrow{\text{vis}} T_1$

Replicated Data Bases

Consistency Models

Φ	Consistency model	Axioms (Figure 2)	Fractured reads	Causality violation	Lost update	Long fork	Write skew	
RA	Read Atomic [6]	INT, EXT	✗	✓	✓	✓	✓	RA
CC	Causal consistency [12, 19]	INT, EXT, TRANSVIS	✗	✗	✓	✓	✓	CC
PSI	Parallel snapshot isolation [21, 24]	INT, EXT, TRANSVIS, NoCONFLICT	✗	✗	✗	✓	✓	PSI
PC	Prefix consistency [13]	INT, EXT, PREFIX	✗	✗	✓	✗	✓	PC
SI	Snapshot isolation [8]	INT, EXT, PREFIX, NoCONFLICT	✗	✗	✗	✗	✓	SI
SER	Serialisability [20]	INT, EXT, TOTALVIS	✗	✗	✗	✗	✗	SER

```

graph TD
    RA[RA] --- CC[CC]
    CC --- PC[PC]
    CC --- PSI[PSI]
    PC --- SI[SI]
    SI --- SER[SER]
  
```

Replicated Data Bases

Consistency Models

Φ	Consistency model	Axioms (Figure 2)	Fractured reads	Causality violation	Lost update	Long fork	Write skew	
RA	Read Atomic [6]	INT, EXT	✗	✓	✓	✓	✓	RA
CC	Causal consistency [12, 19]	INT, EXT, TRANSVIS	✗	✗	✓	✓	✓	CC
PSI	Parallel snapshot isolation [21, 24]	INT, EXT, TRANSVIS, NoCONFLICT	✗	✗	✗	✓	✓	PSI
PC	Prefix consistency [13]	INT, EXT, PREFIX	✗	✗	✓	✗	✓	PC
SI	Snapshot isolation [8]	INT, EXT, PREFIX, NoCONFLICT	✗	✗	✗	✗	✓	SI
SER	Serialisability [20]	INT, EXT, TOTALVIS	✗	✗	✗	✗	✗	SER

$$CC \equiv INT \wedge EXT \wedge SER_{TOTAL}$$

$$PSI \equiv CC \wedge CONFLICT$$

$$PC \equiv CC \wedge PREFIX$$

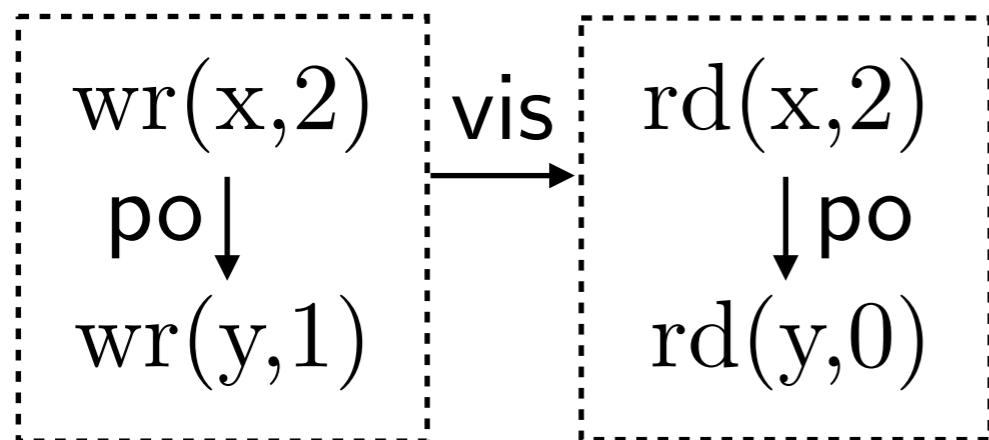
$$SI \equiv PSI \wedge PREFIX$$

$$SER \equiv INT \wedge EXT \wedge TOTAL_{HB}$$

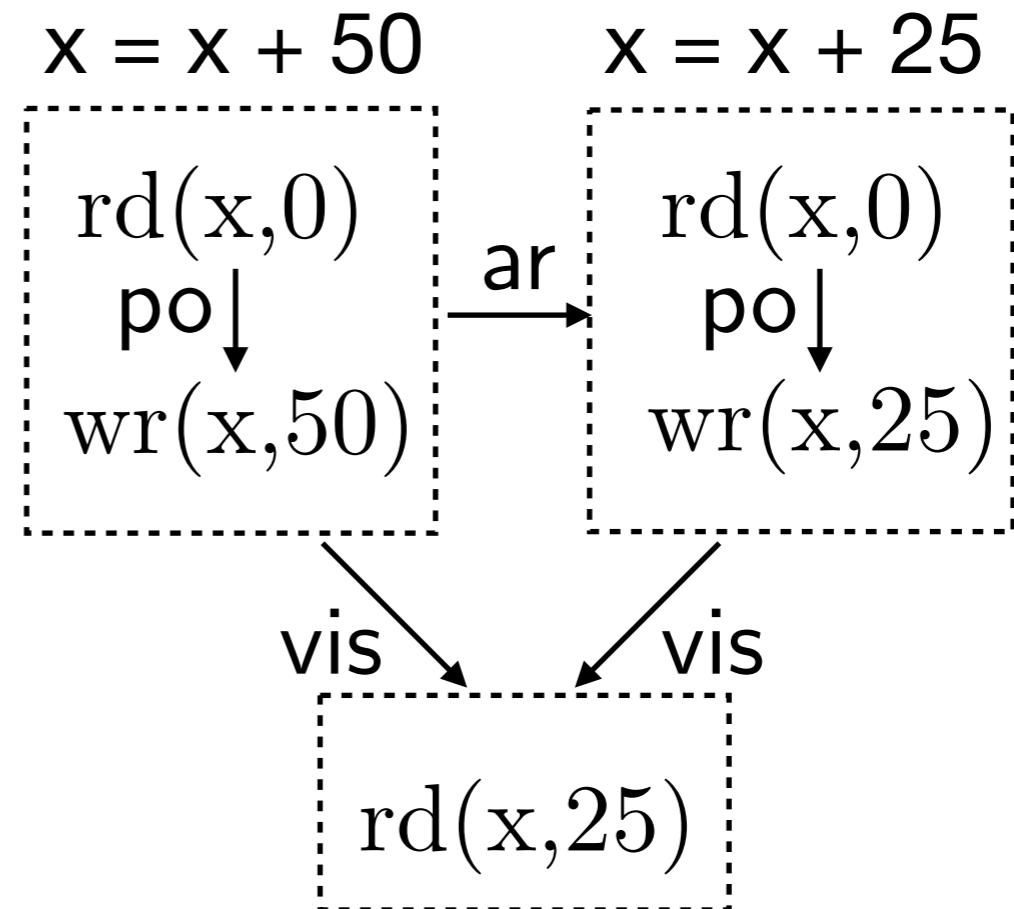
Replicated Data Bases

Anomalies

Fractured Reads



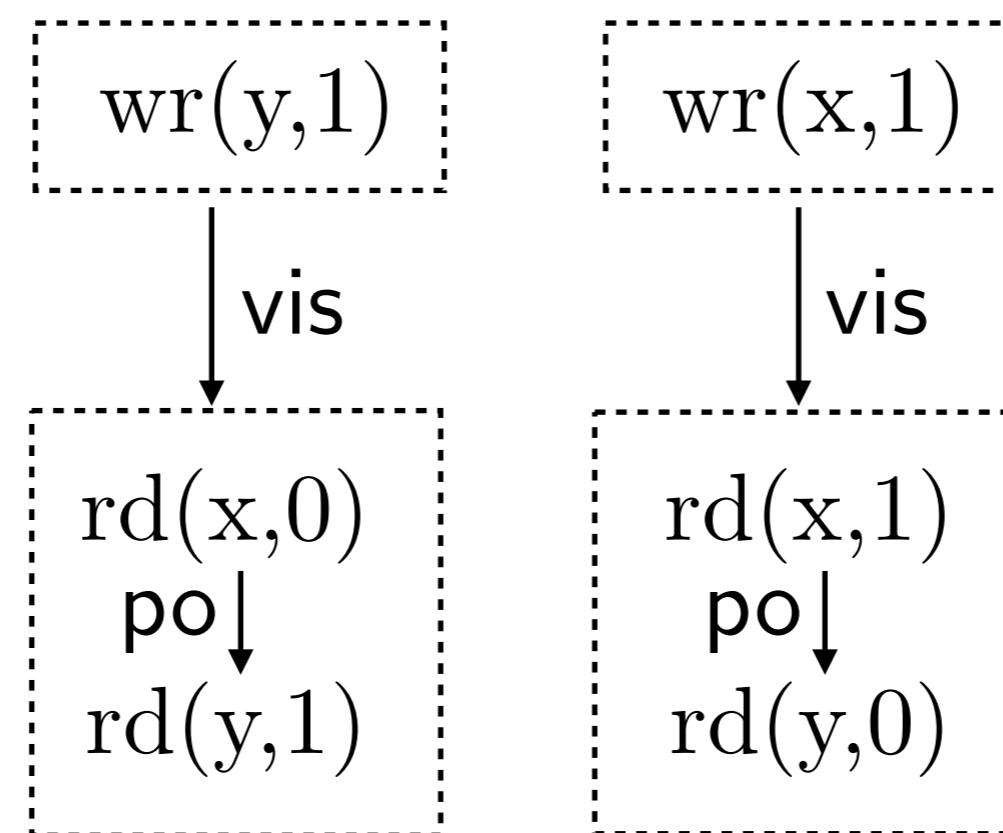
Lost Update



Replicated Data Bases

Anomalies

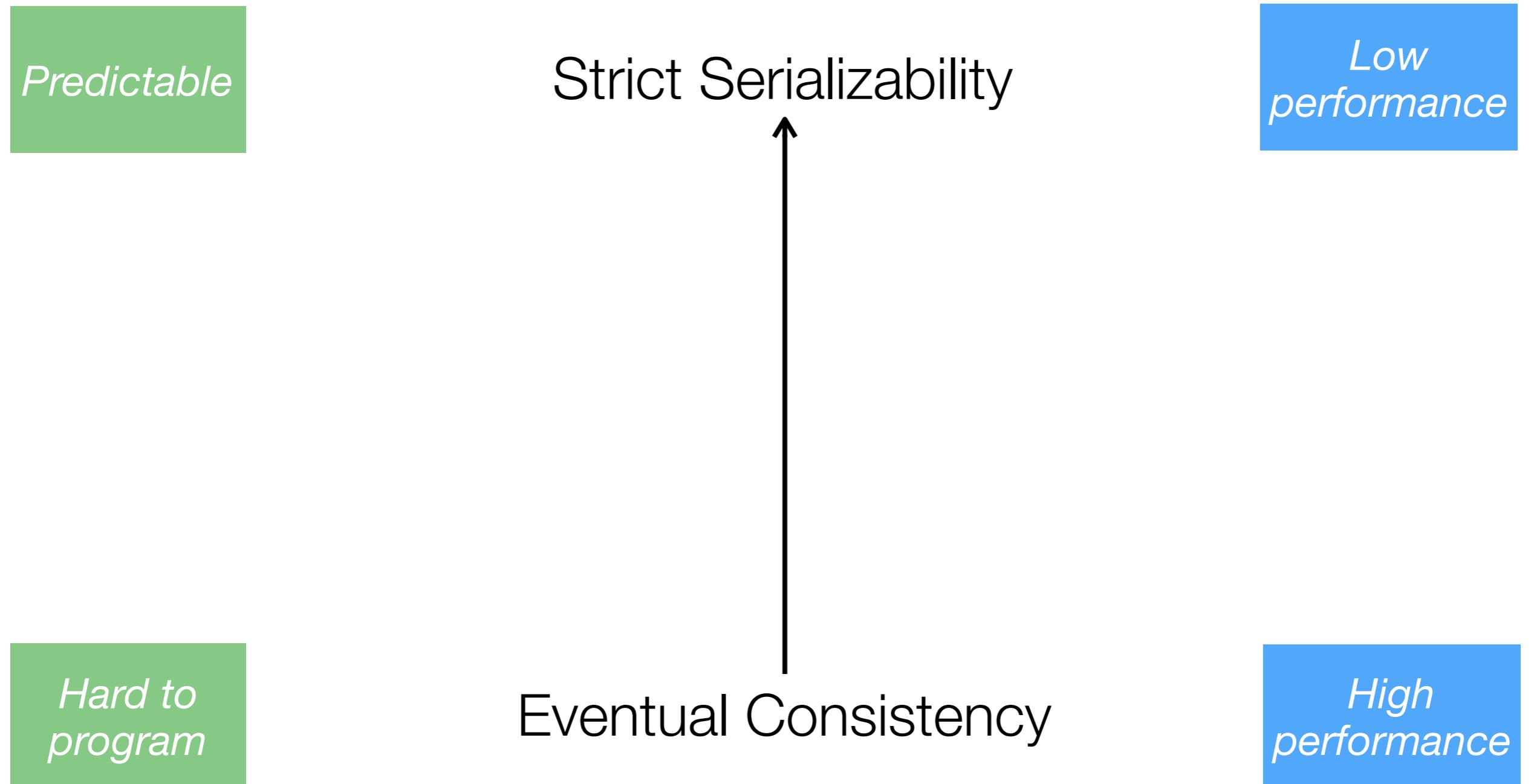
Long Fork (IRIW)



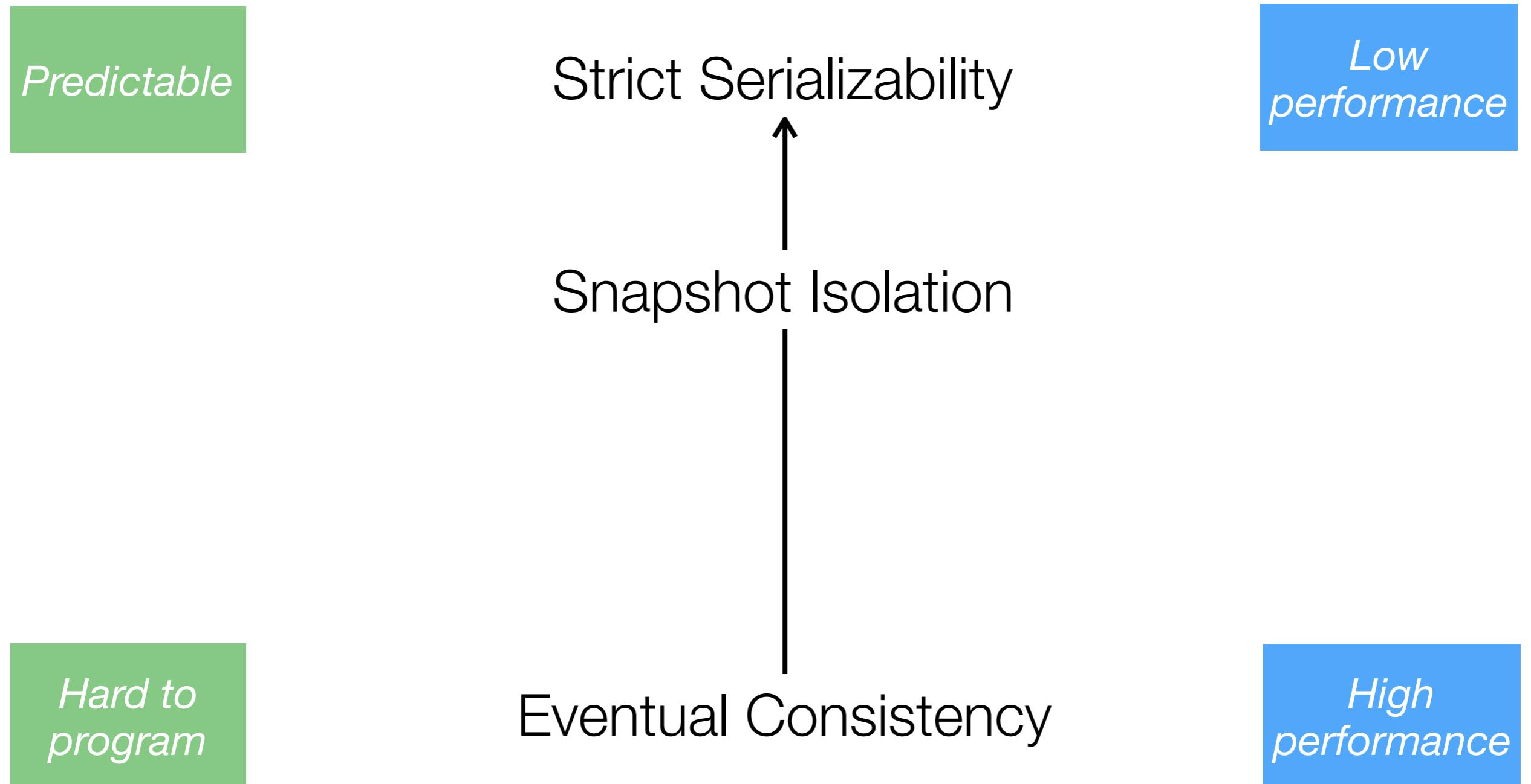
Tiny Demo

A Forest of Models

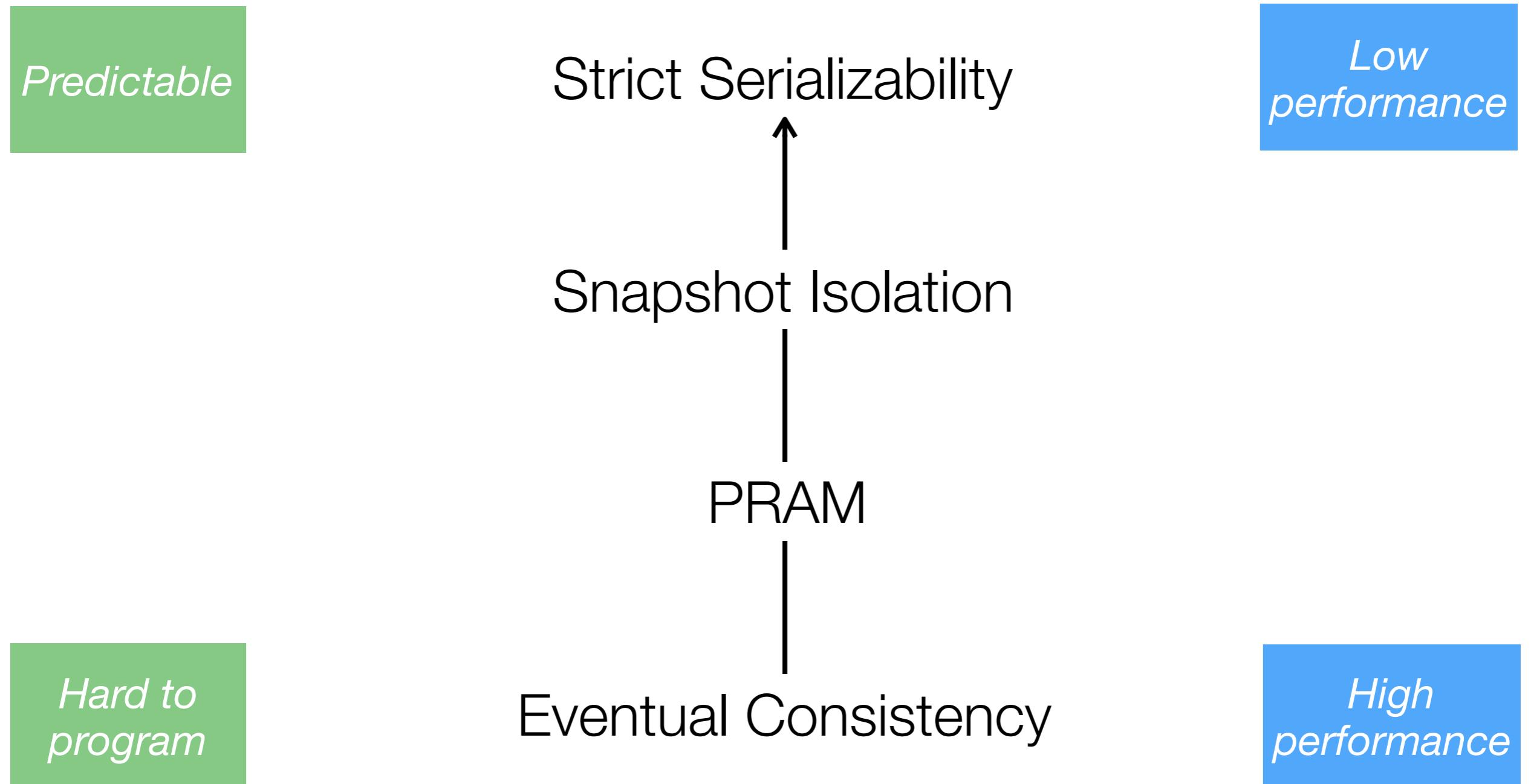
Strong vs. weak?



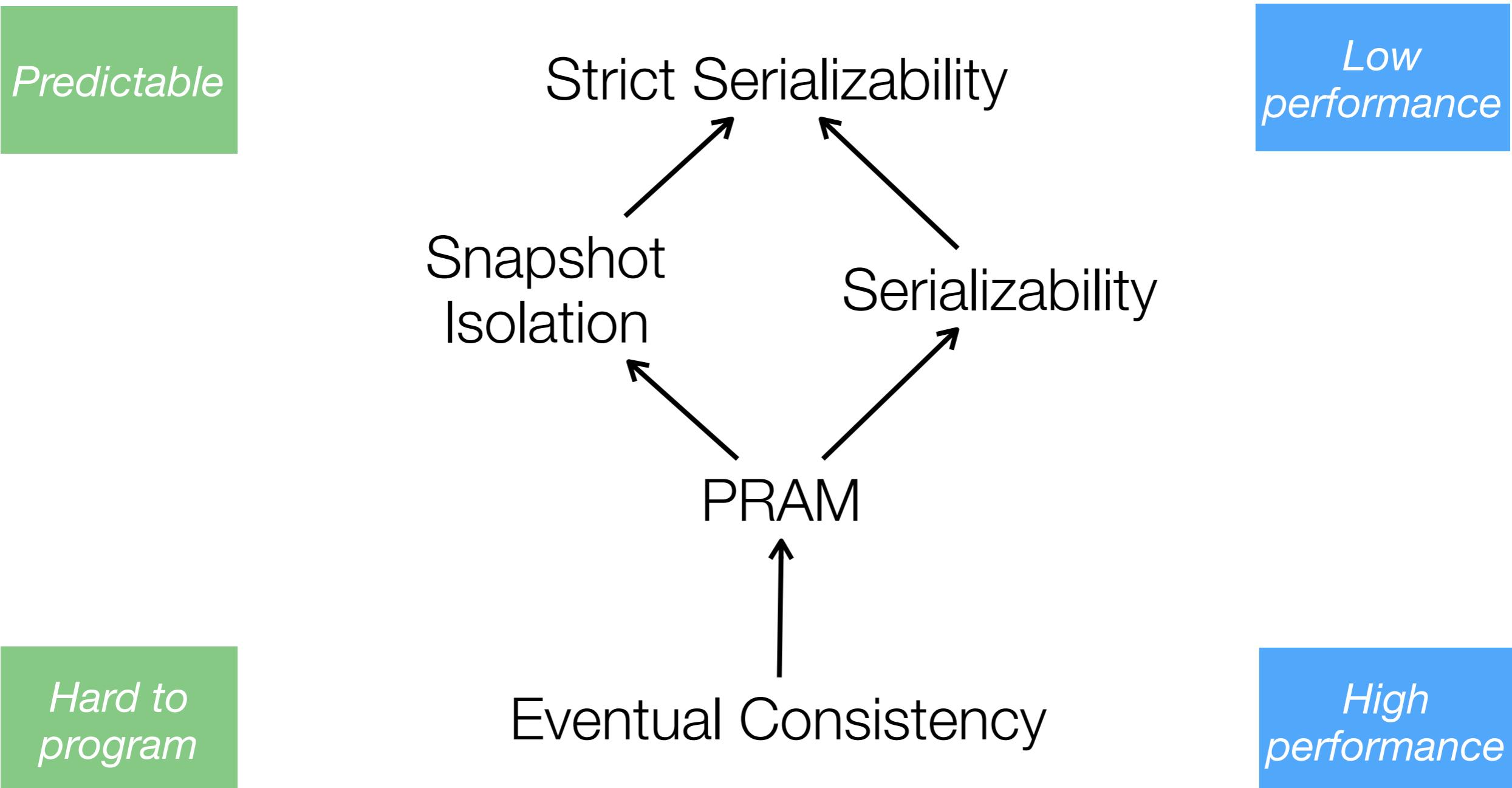
Strong vs. weak?



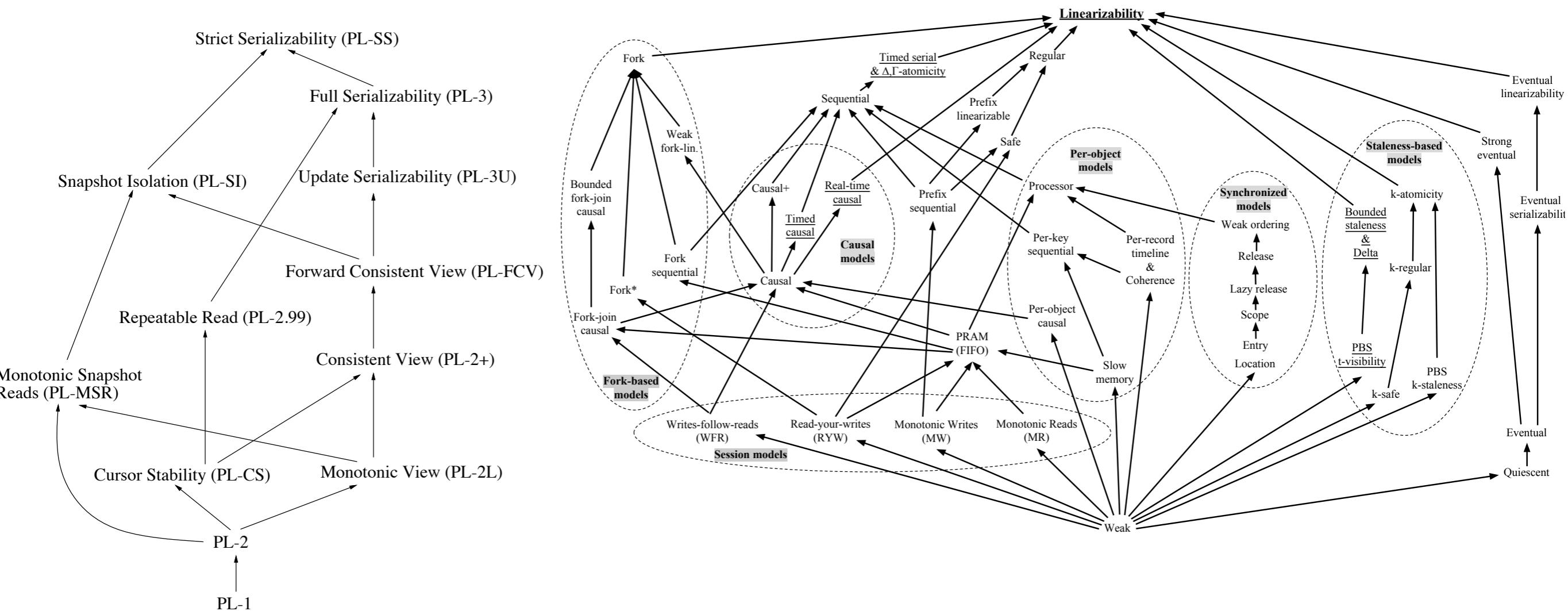
Strong vs. weak?



Strong vs. weak?



Strong vs. weak?



Transactional
Adya 1999

Non-transactional
Viotti & Vukolić 2016

Consistency & Invariants

Consistency & Invariants

Consistency in 3D*

Marc Shapiro¹, Masoud Saeida Ardekani², and Gustavo Petri³

¹ Sorbonne-Universités-UPMC-LIP6 & Inria Paris

² Purdue University [†]

³ IRIF, Université Paris Diderot

Abstract

Comparisons of different consistency models often try to place them in a linear strong-to-weak order. However this view is clearly inadequate, since it is well known, for instance, that Snapshot Isolation and Serialisability are incomparable. In the interest of a better understanding, we propose a new classification, along three dimensions, related to: a total order of writes, a causal order of reads, and transactional composition of multiple operations. A model may be stronger than another on one dimension and weaker on another. We believe that this new classification scheme is both scientifically sound and has good explicative value. The current paper presents the three-dimensional design space intuitively.

1998 ACM Subject Classification C.2.4 Distributed databases; D.1.3 Concurrent programming;
D.2.4 Software/Program Verification; E.1 Distributed data structures

Keywords and phrases Consistency models; Replicated data; Structural invariants; Correctness of distributed systems;

Digital Object Identifier 10.4230/LIPIcs.CONCUR.2016.<article-no>

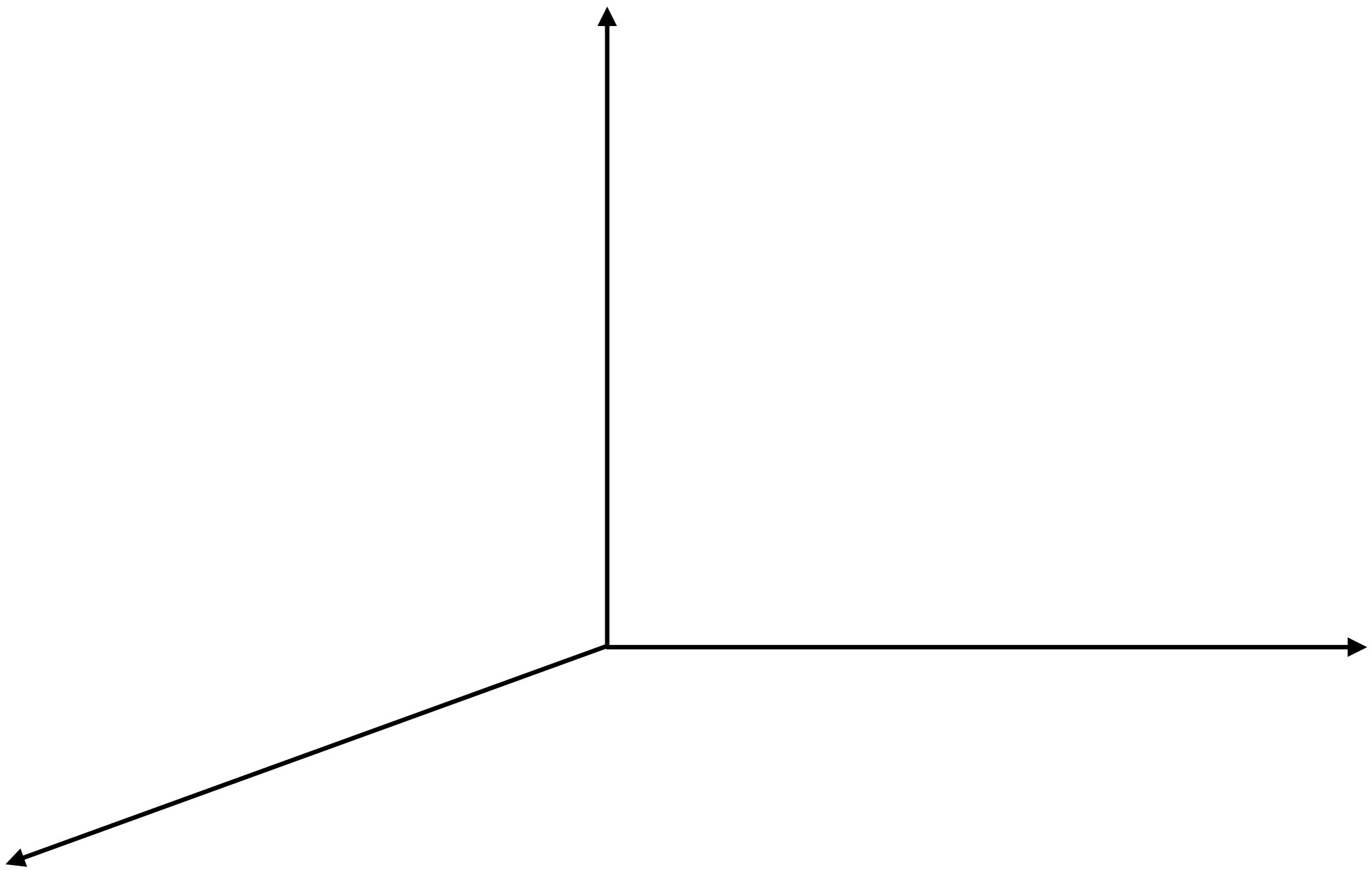
Consistency & Invariants

- Consistency in 3D
 - Characterization of consistency models according to the guarantees they provide
- Dimensions of Guarantees
 - Single object
 - Propagation of effects on different objects
 - Composition of objects

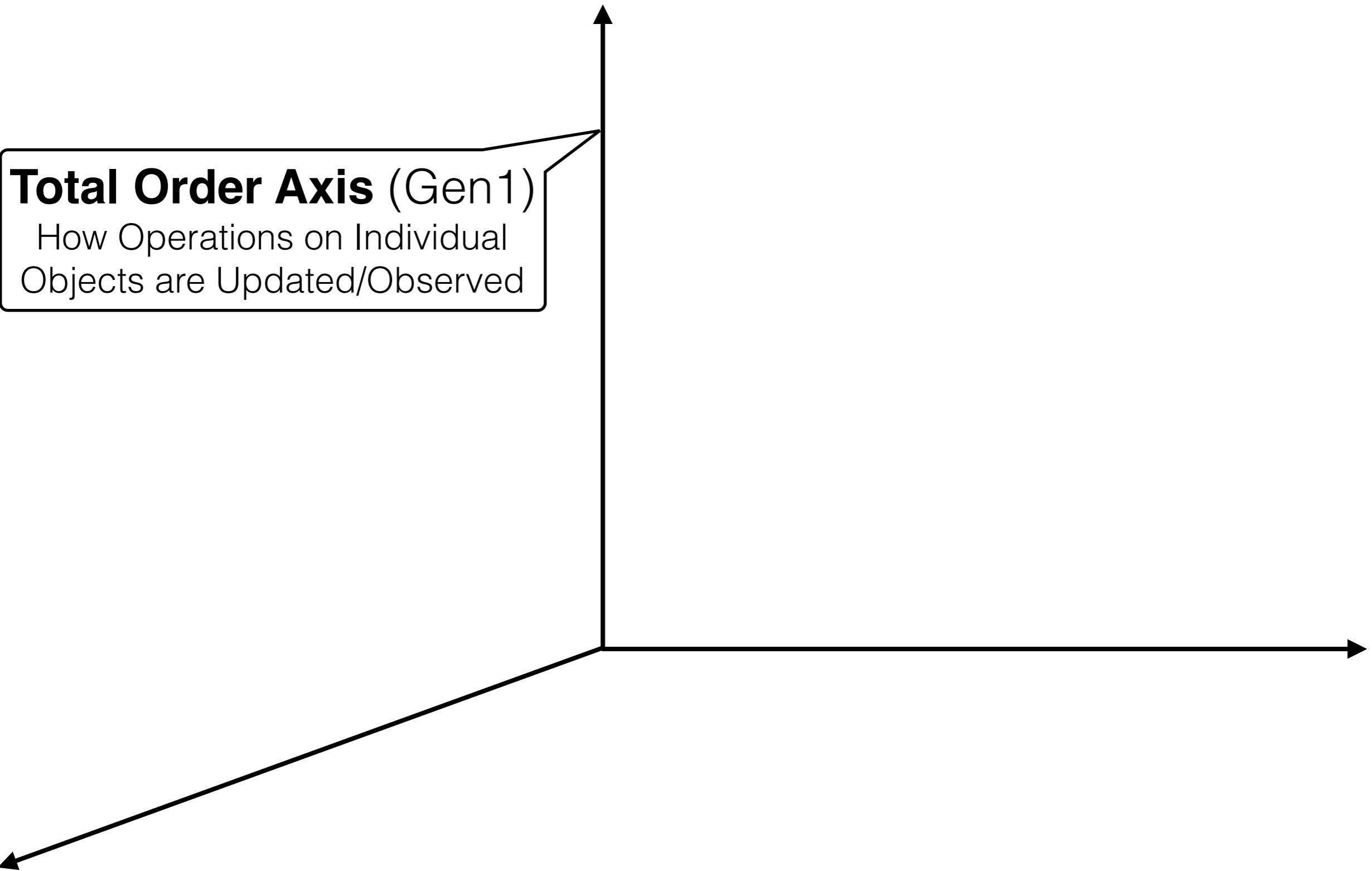
Three classes...

	...of invariant	... of protocol
Gen1	Constrain value of an object	Total order of operations
PO	Ordering between operations	Visibility
EQ	State equivalence between objects	Composition

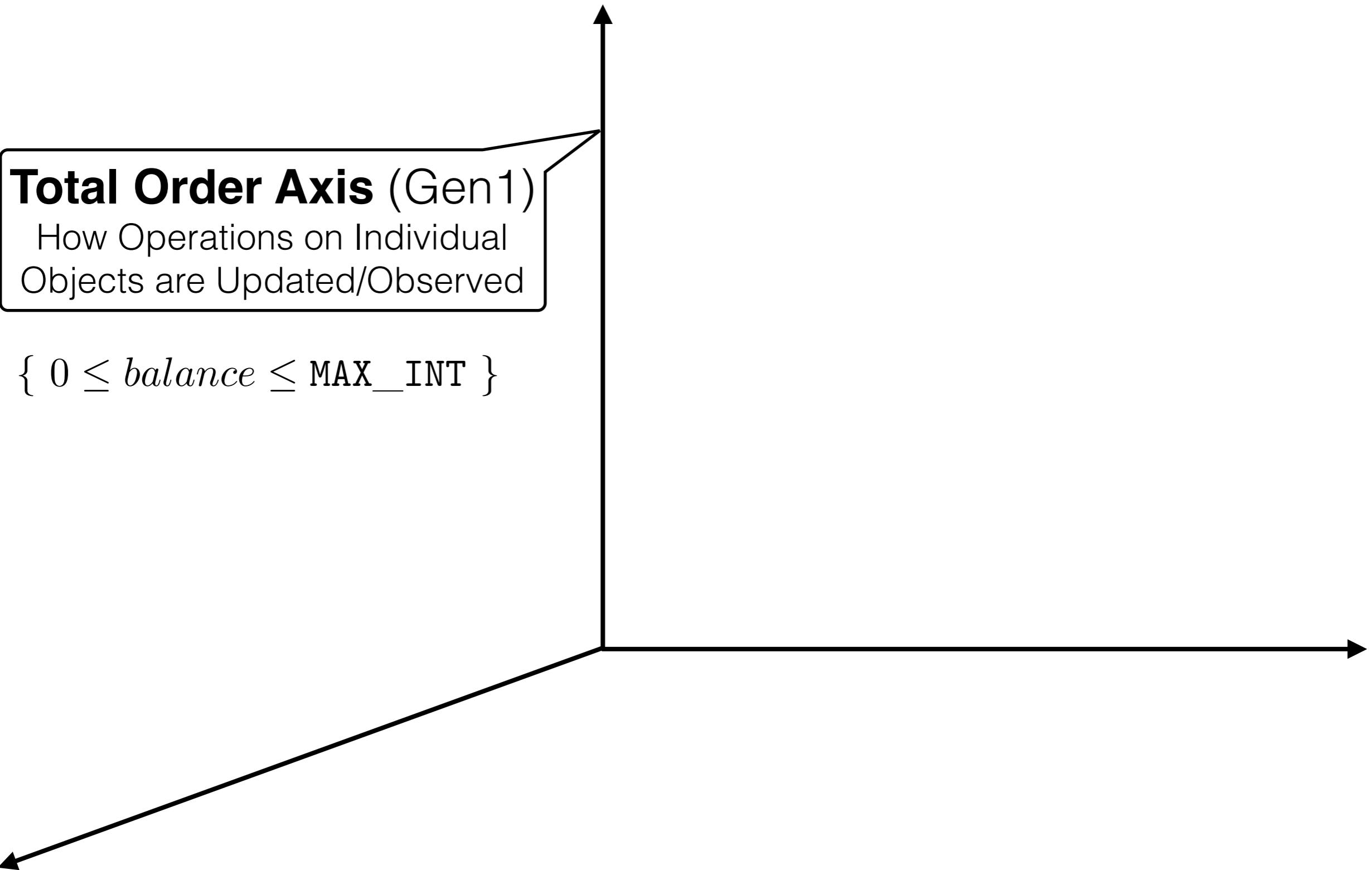
Consistency in 3D



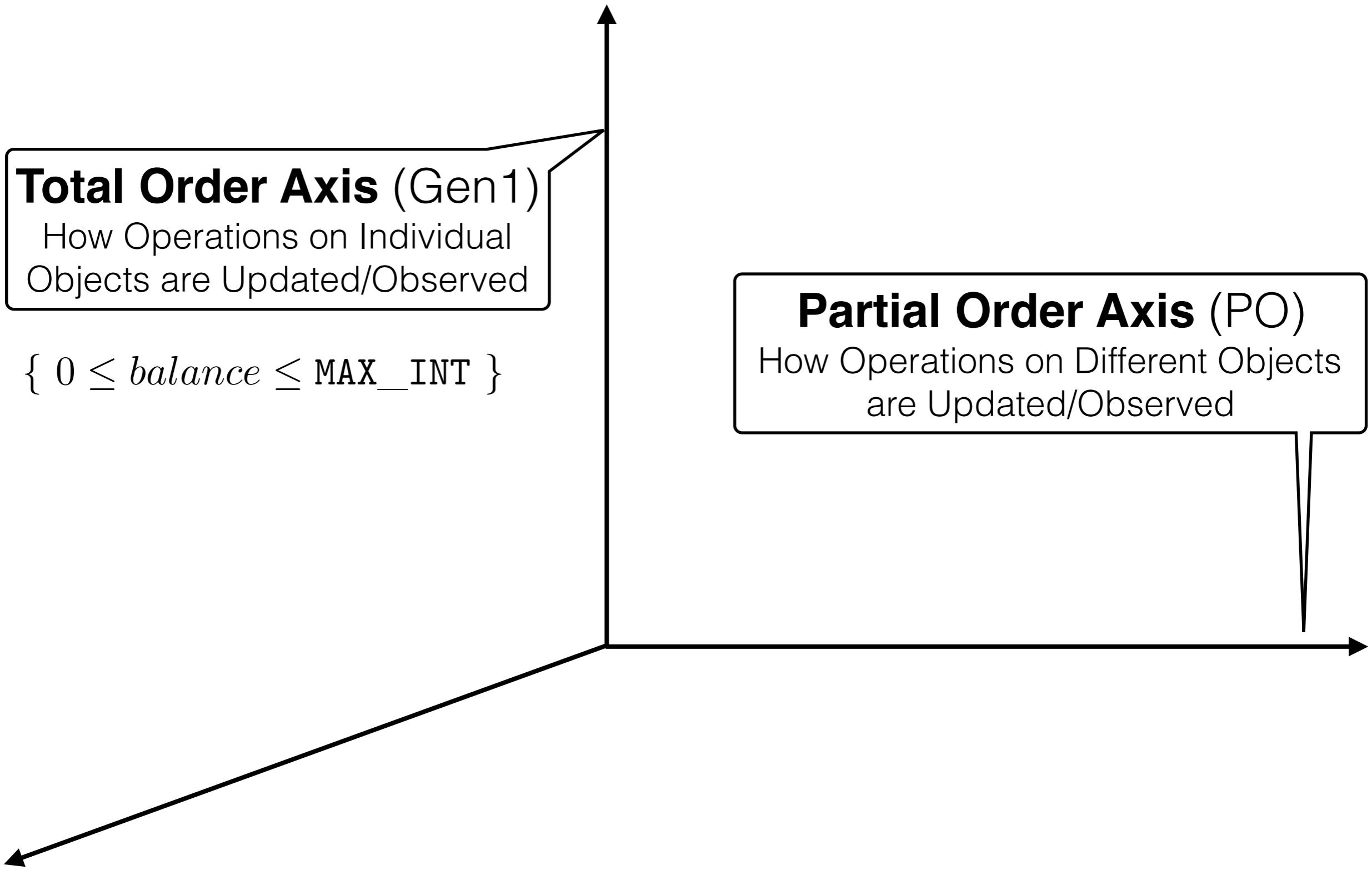
Consistency in 3D



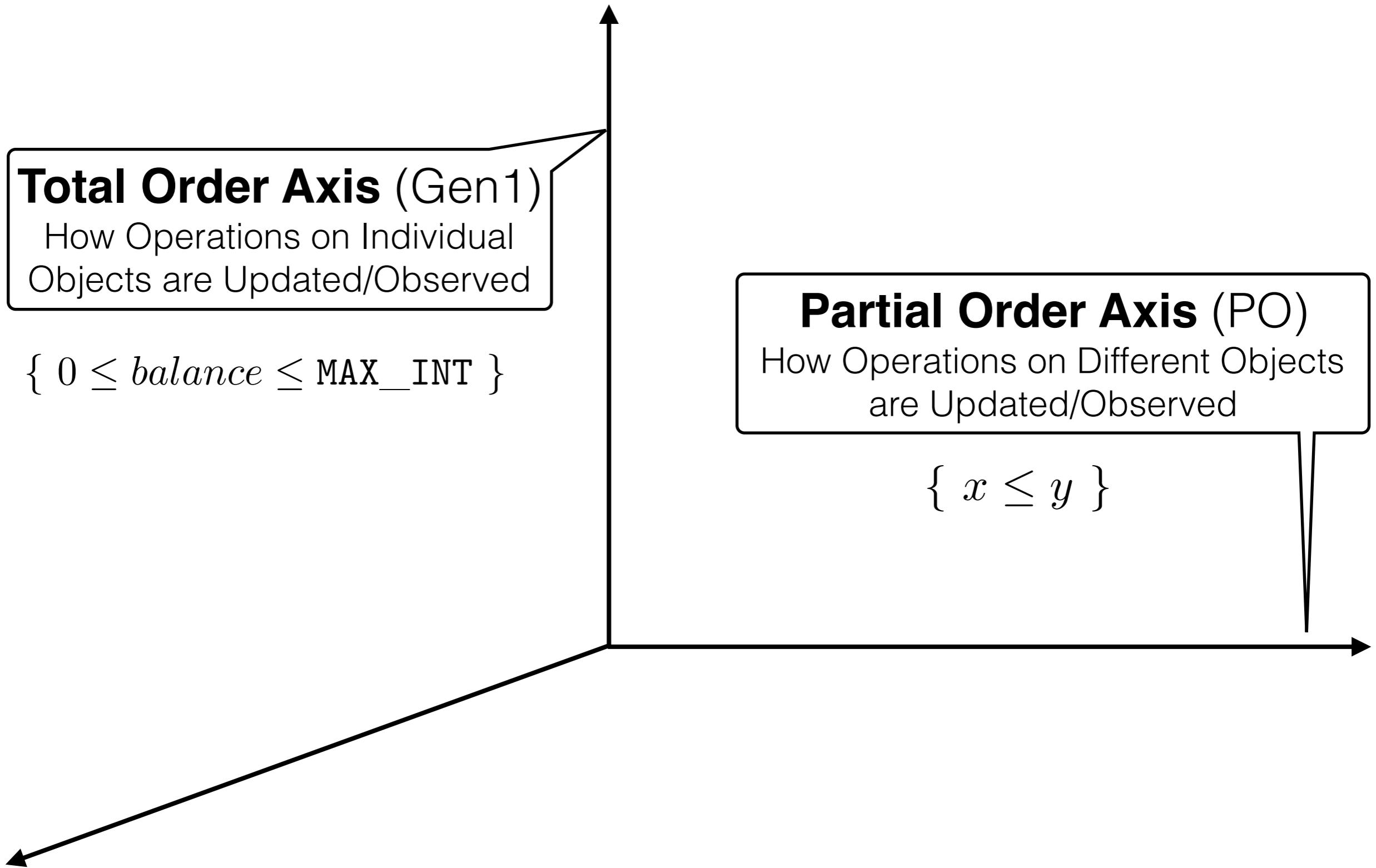
Consistency in 3D



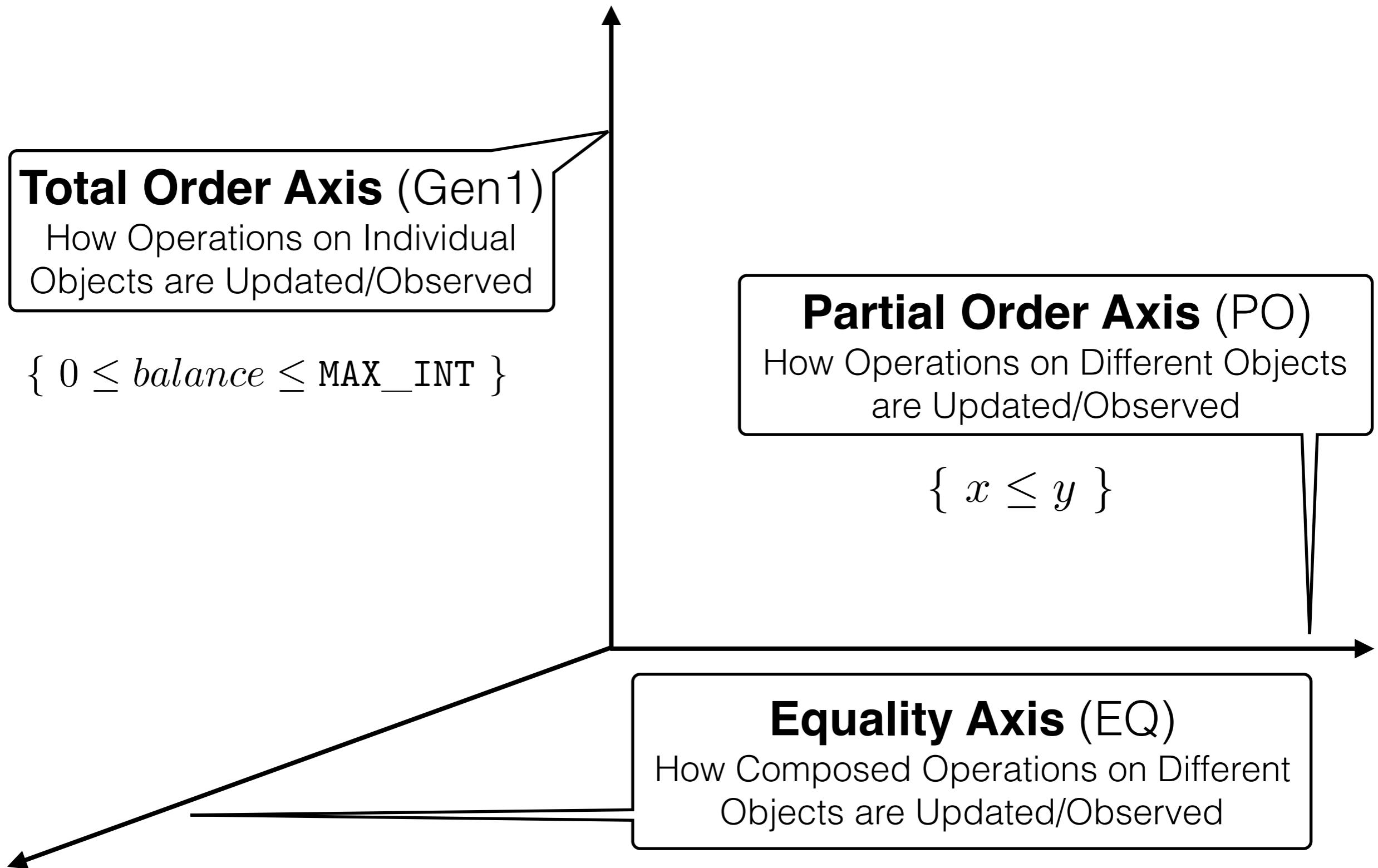
Consistency in 3D



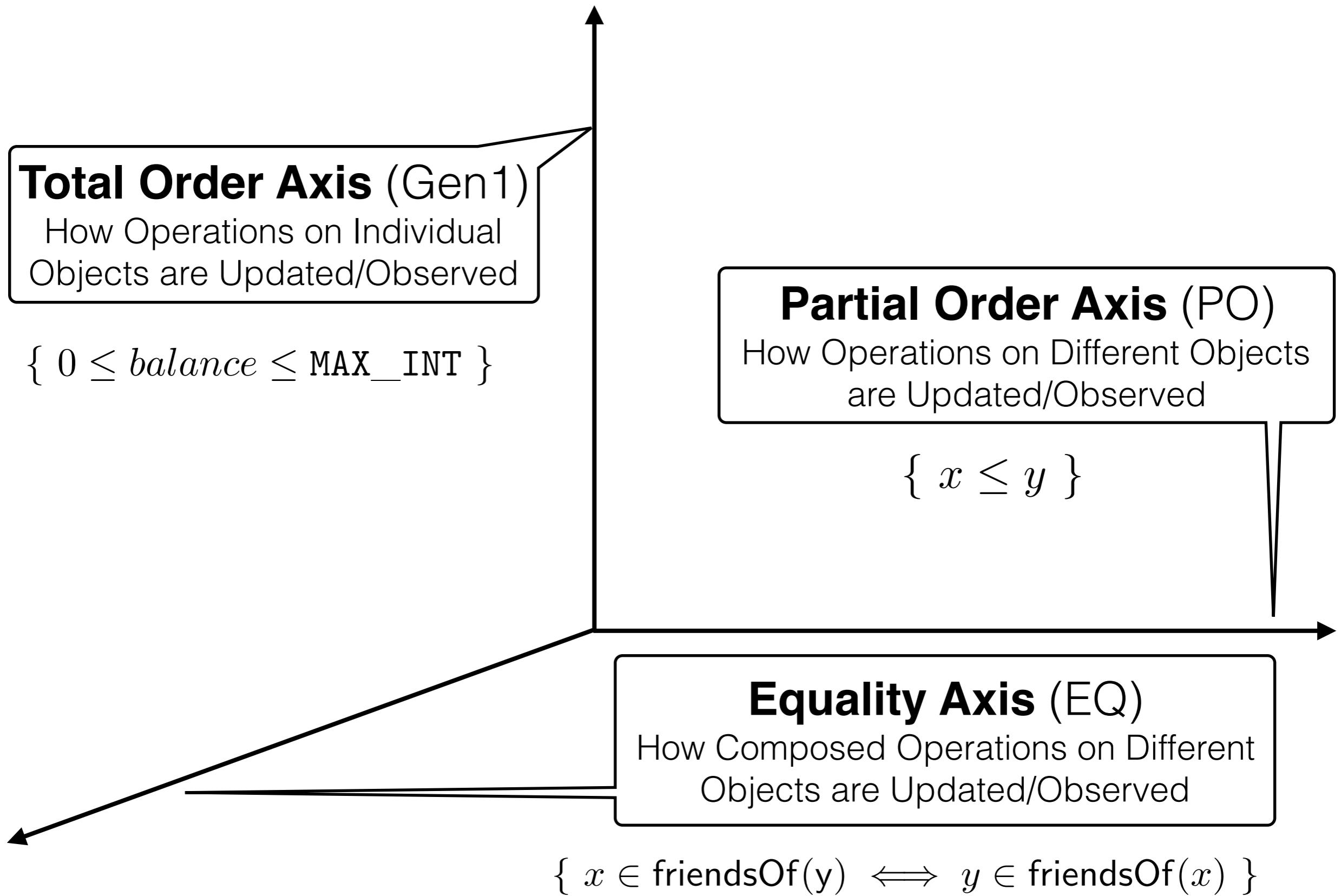
Consistency in 3D



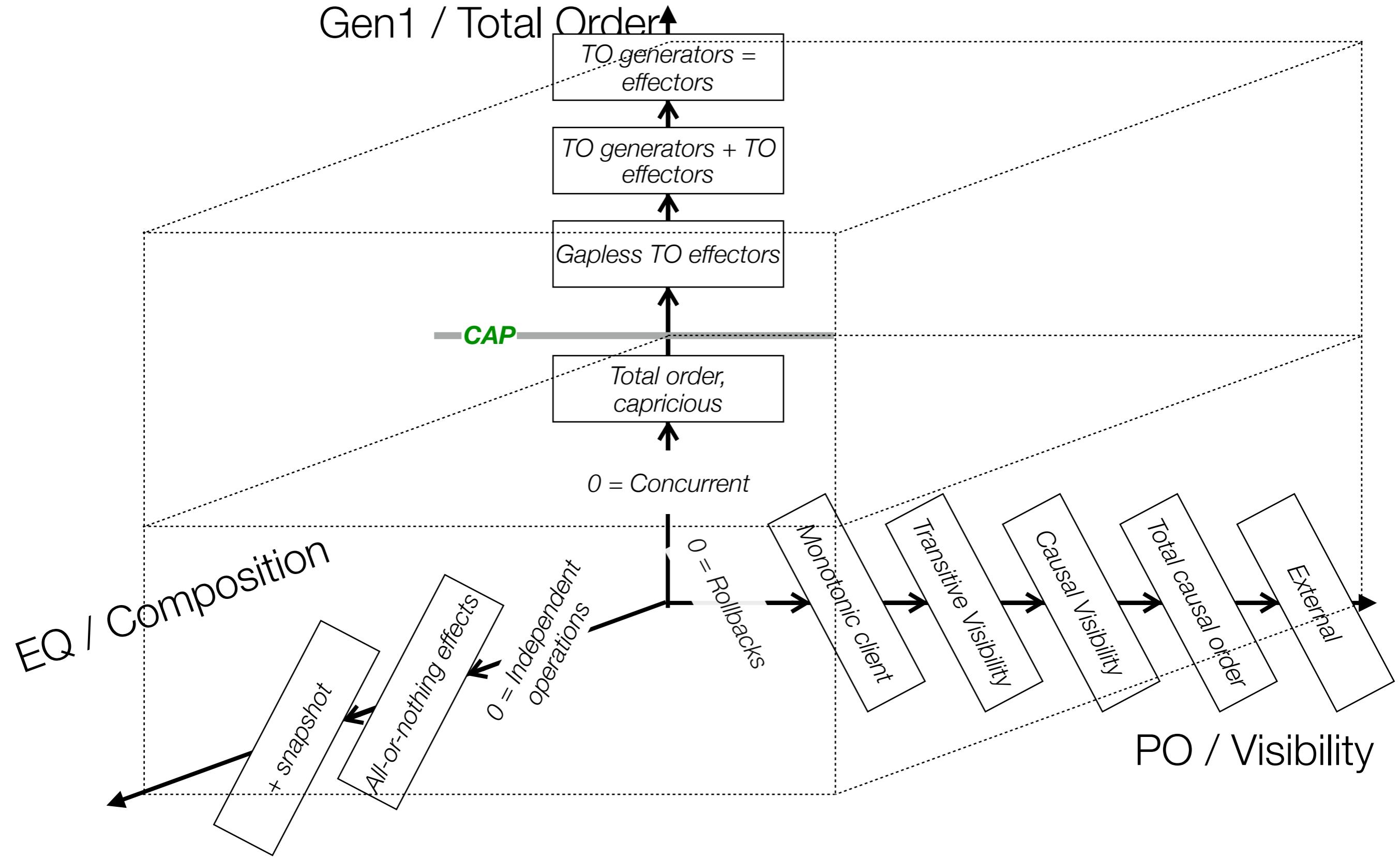
Consistency in 3D



Consistency in 3D

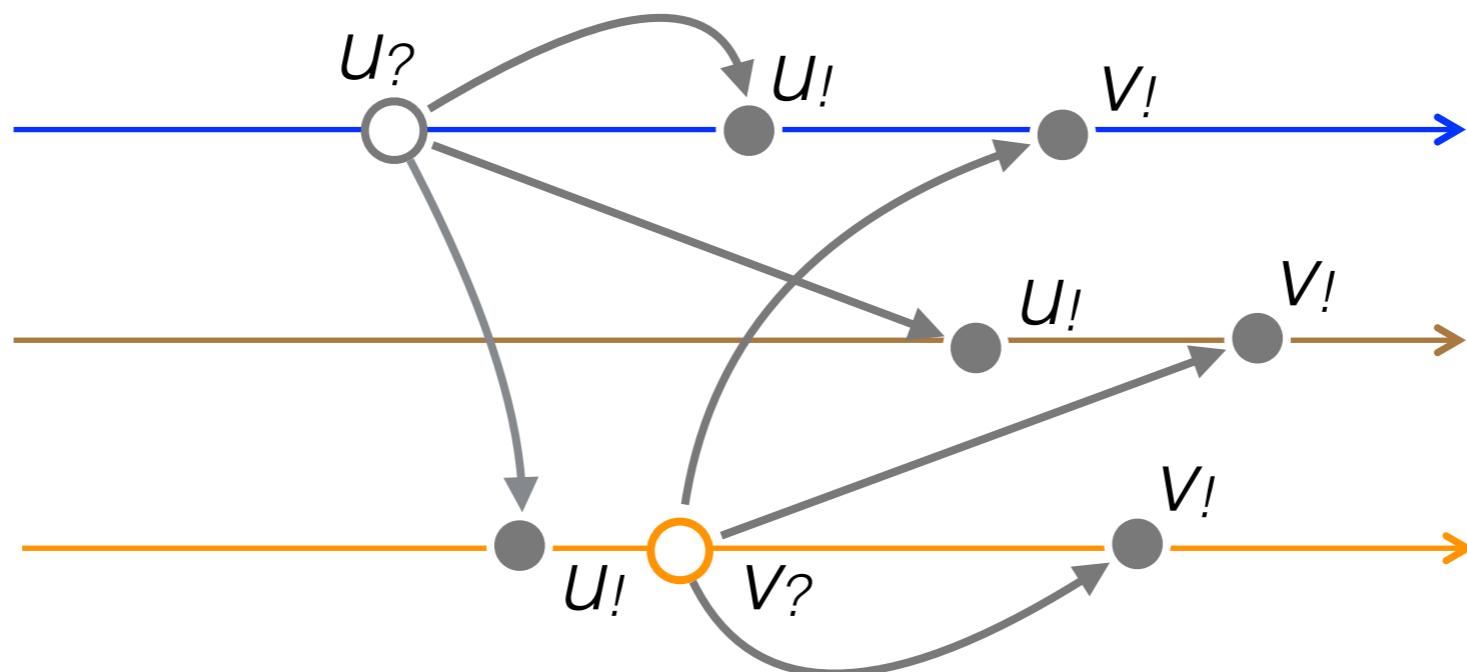


Three dimensions



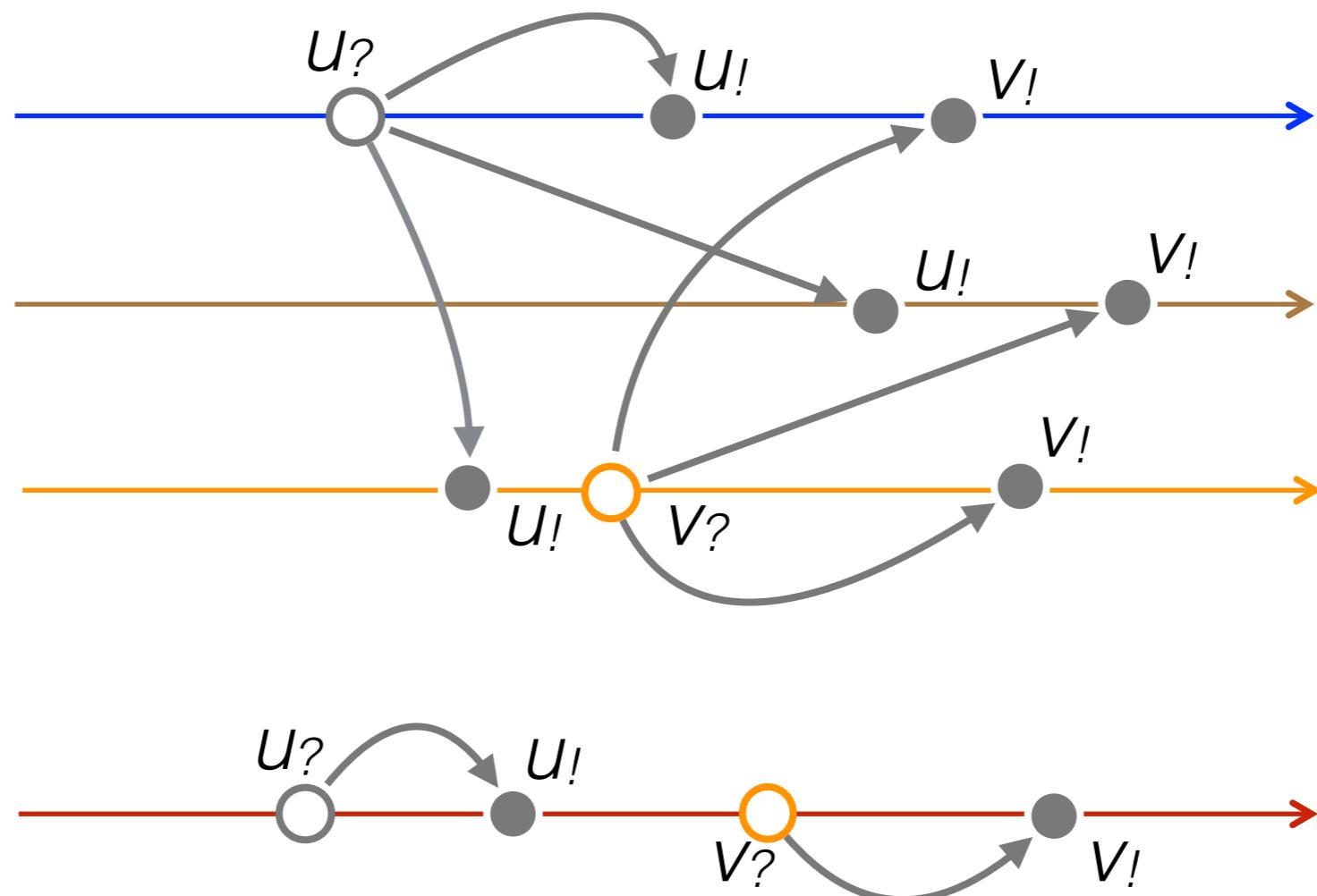
Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE=TOG)



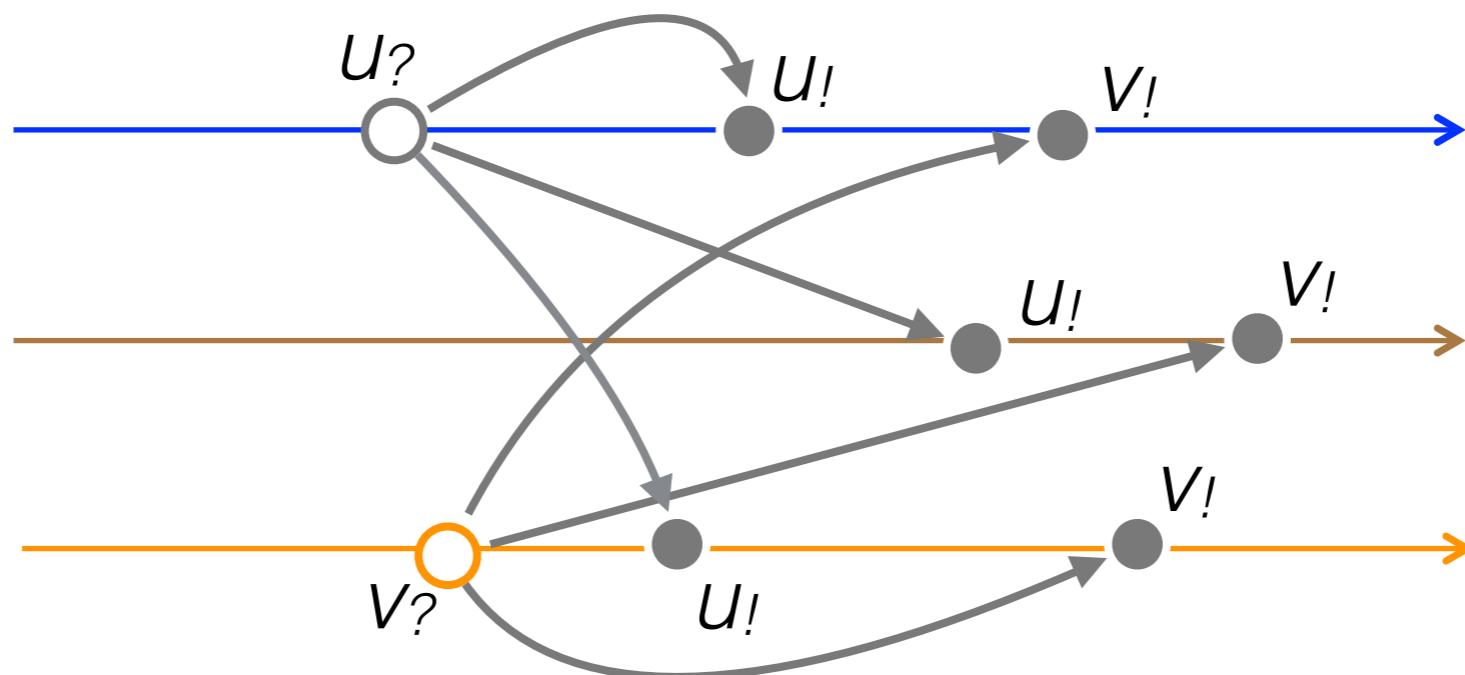
Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE=TOG)



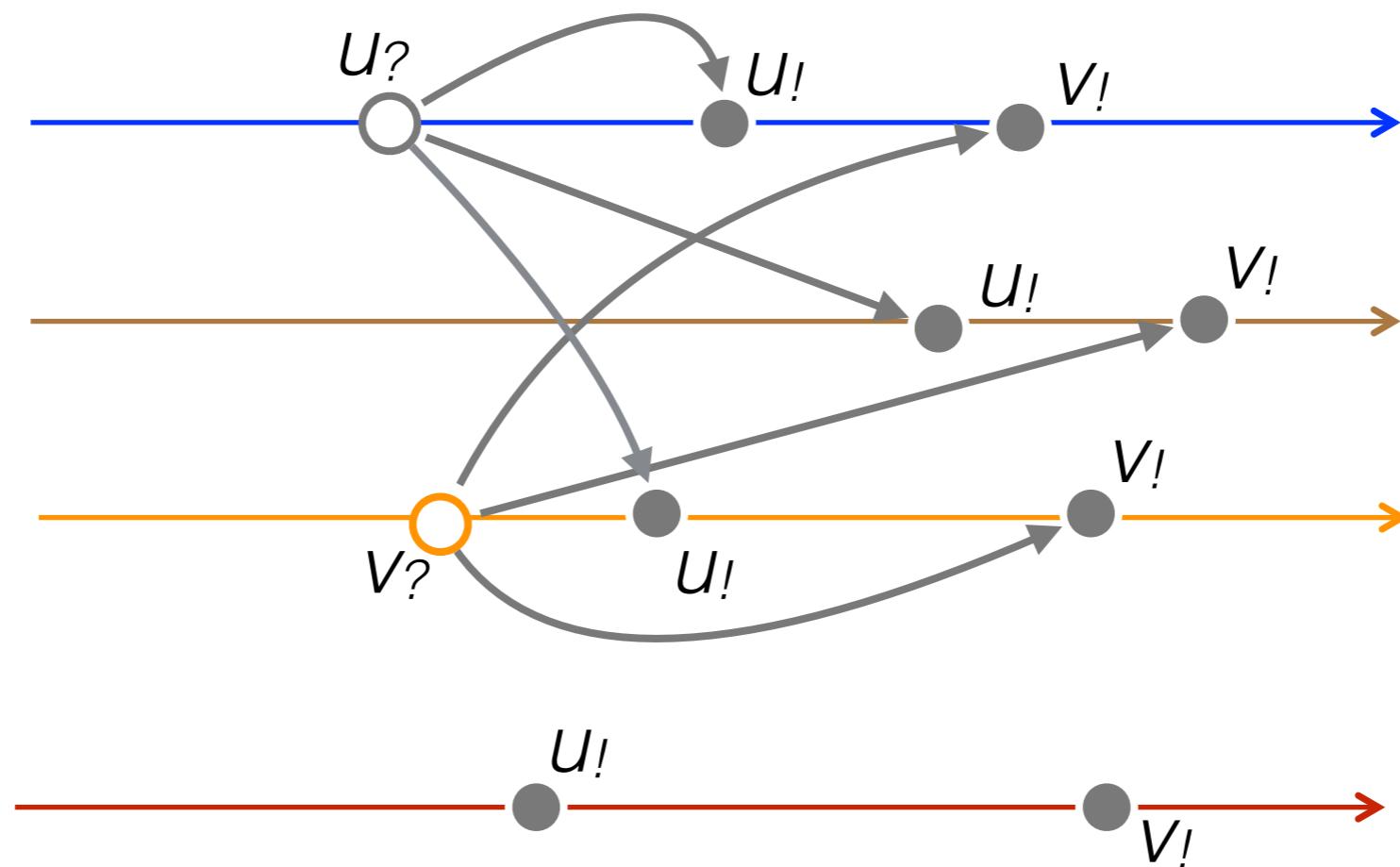
Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE_1)



Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE_1)



Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE₁)

Total Order Axis

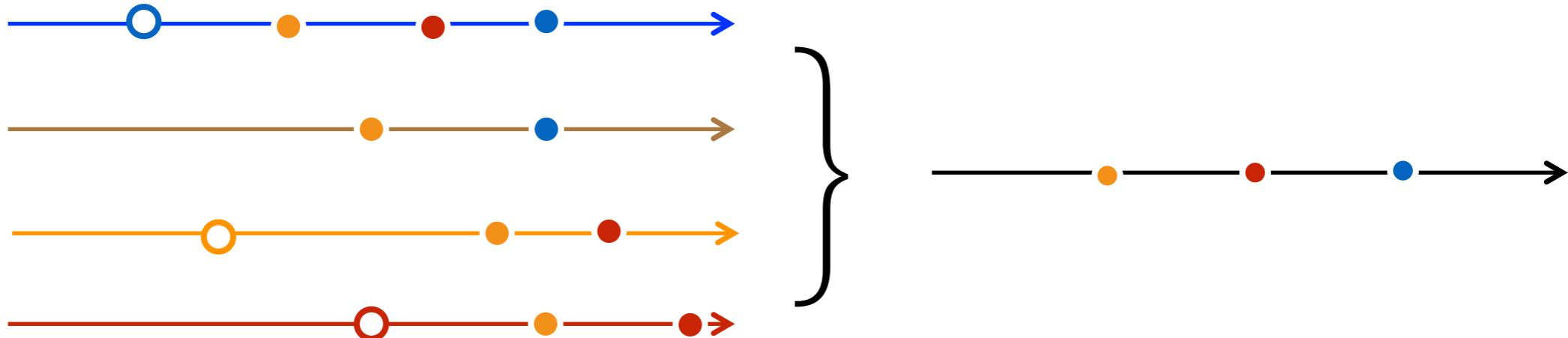
- Assumption: Single Object
- Total Order of Effectors and Generators (TOE₁)
 - Gapless TOE₁: all replicas apply all effectors in the same order

Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE₁)
 - Gapless TOE₁: all replicas apply all effectors in the same order
 - Capricious TOE₁: replicas apply a subset of the effectors in an order consistent with a global total order

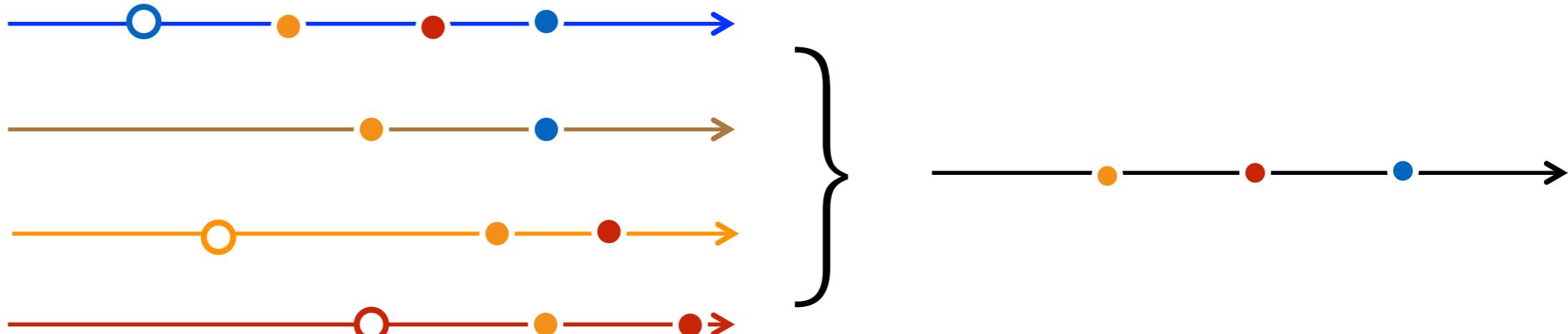
Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE₁)
 - Gapless TOE₁: all replicas apply all effectors in the same order
 - Capricious TOE₁: replicas apply a subset of the effectors in an order consistent with a global total order

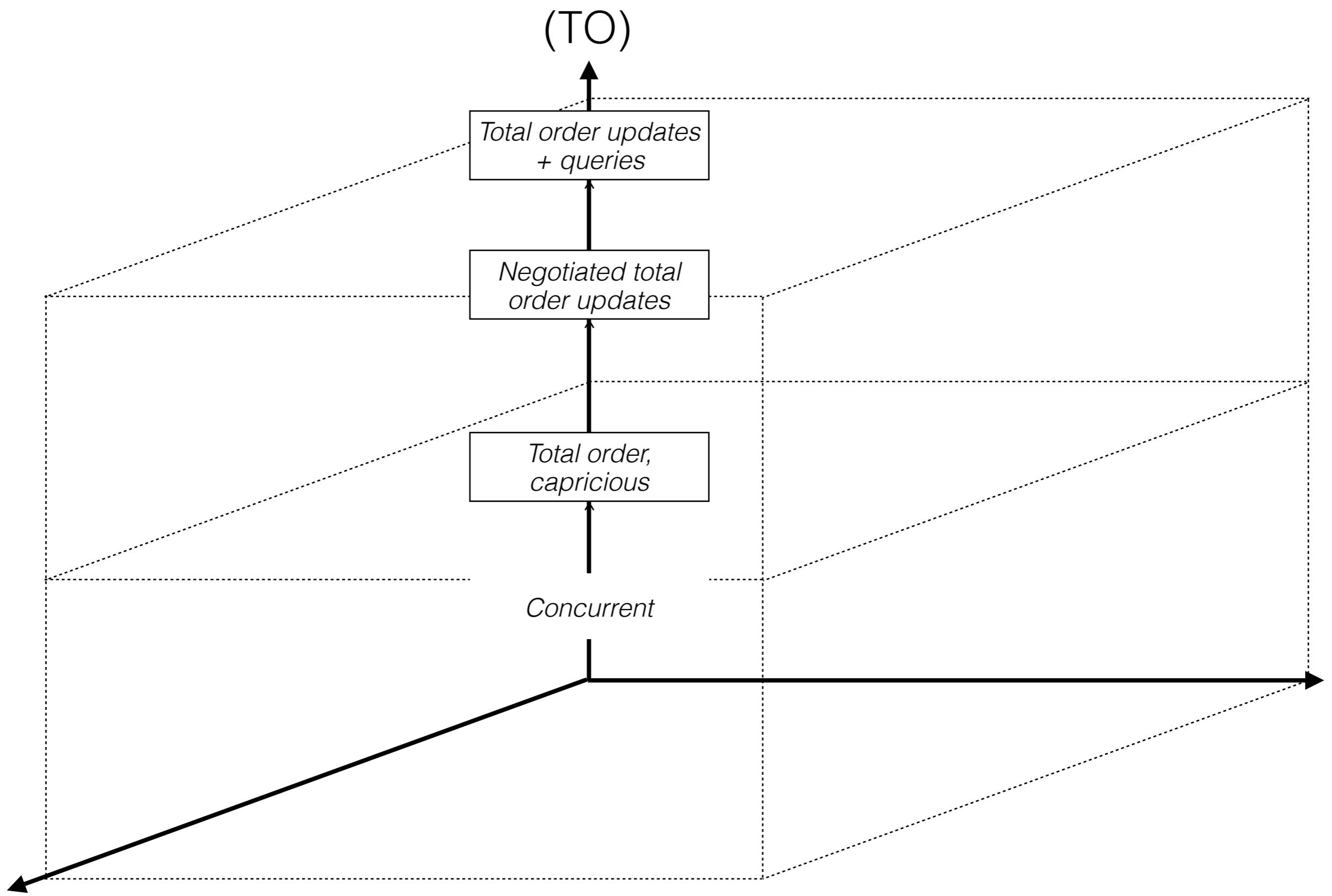


Total Order Axis

- Assumption: Single Object
- Total Order of Effectors and Generators (TOE₁)
 - Gapless TOE₁: all replicas apply all effectors in the same order
 - Capricious TOE₁: replicas apply a subset of the effectors in an order consistent with a global total order

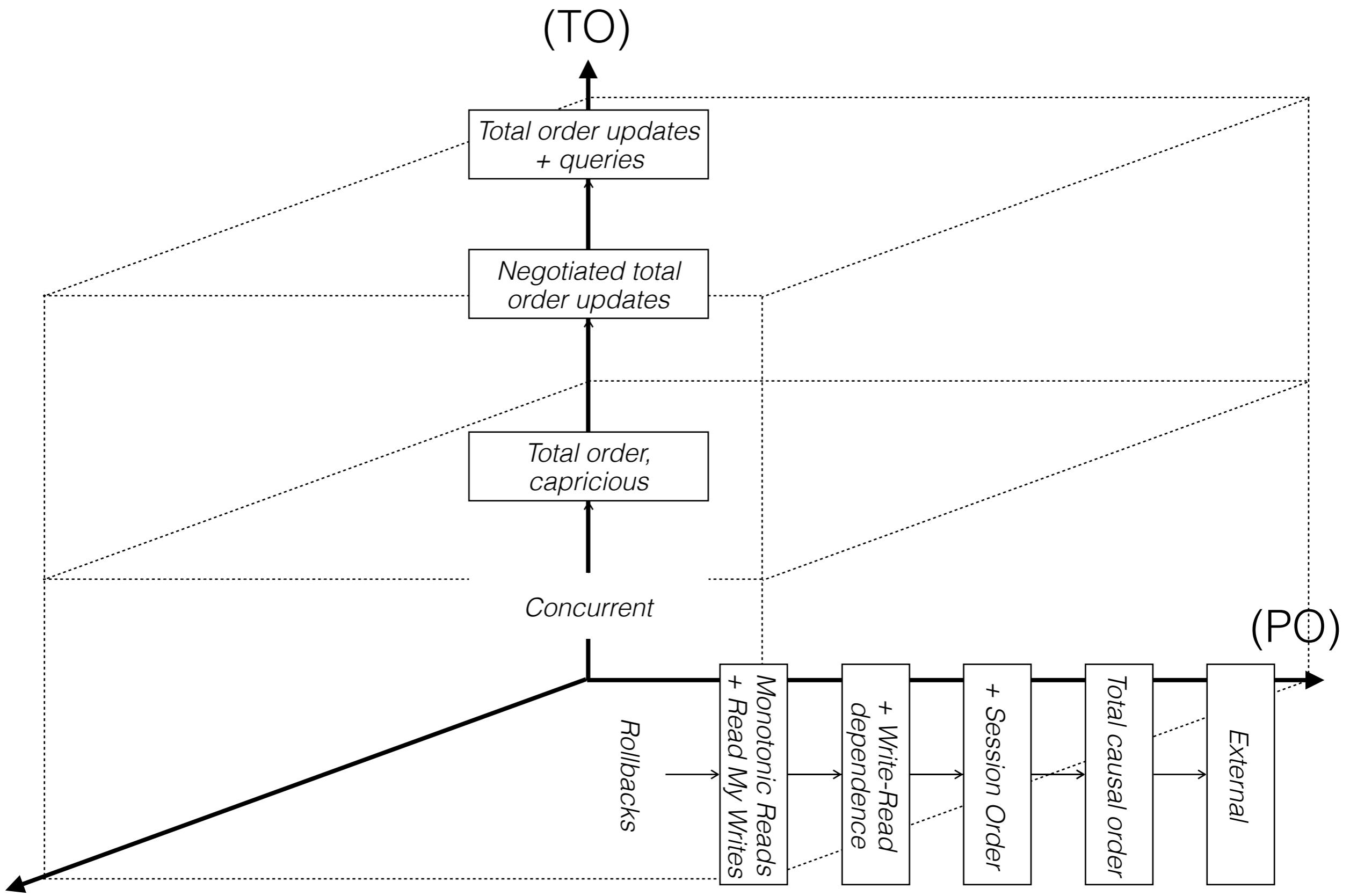


- Concurrent Updates (No Global Ordering)



Partial Order Axis

- Assumption: Multiple (2) Objects
- Client Guarantees:
 - Read Own Writes
 - Monotonicity (Reads/Writes)
 - Preservation of (anti)Dependencies
- Visibility Properties:
 - Transitive Visibility
 - Causal Visibility



Partial Order Axis (Invariants)

- Assumptions:
 - (i) Multiple Object,
 - (ii) State Based,
 - (iii) O is a valid object for I
- Invariants Relating Objects
 - $x \leq y$
 - $P(x) \implies Q(y)$
- Programming:
 - Demarcation Protocol
 - Escrow

Equality Order Axis

- Assumption: Multiple (n) Objects
- Transactions

Equality Order Axis

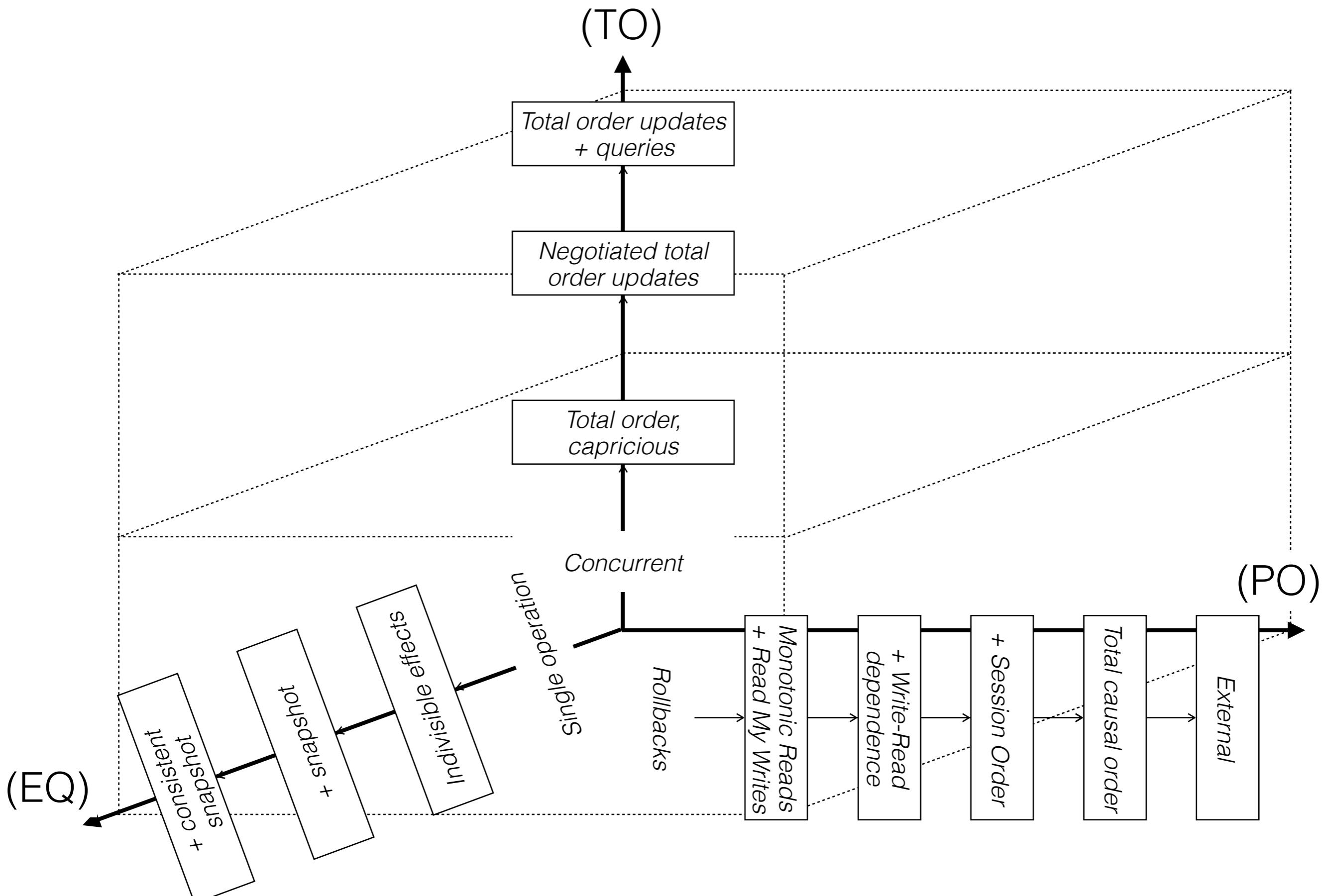
- Assumption: Multiple (n) Objects
- Transactions
 - Write-atomicity: All-or-nothing

Equality Order Axis

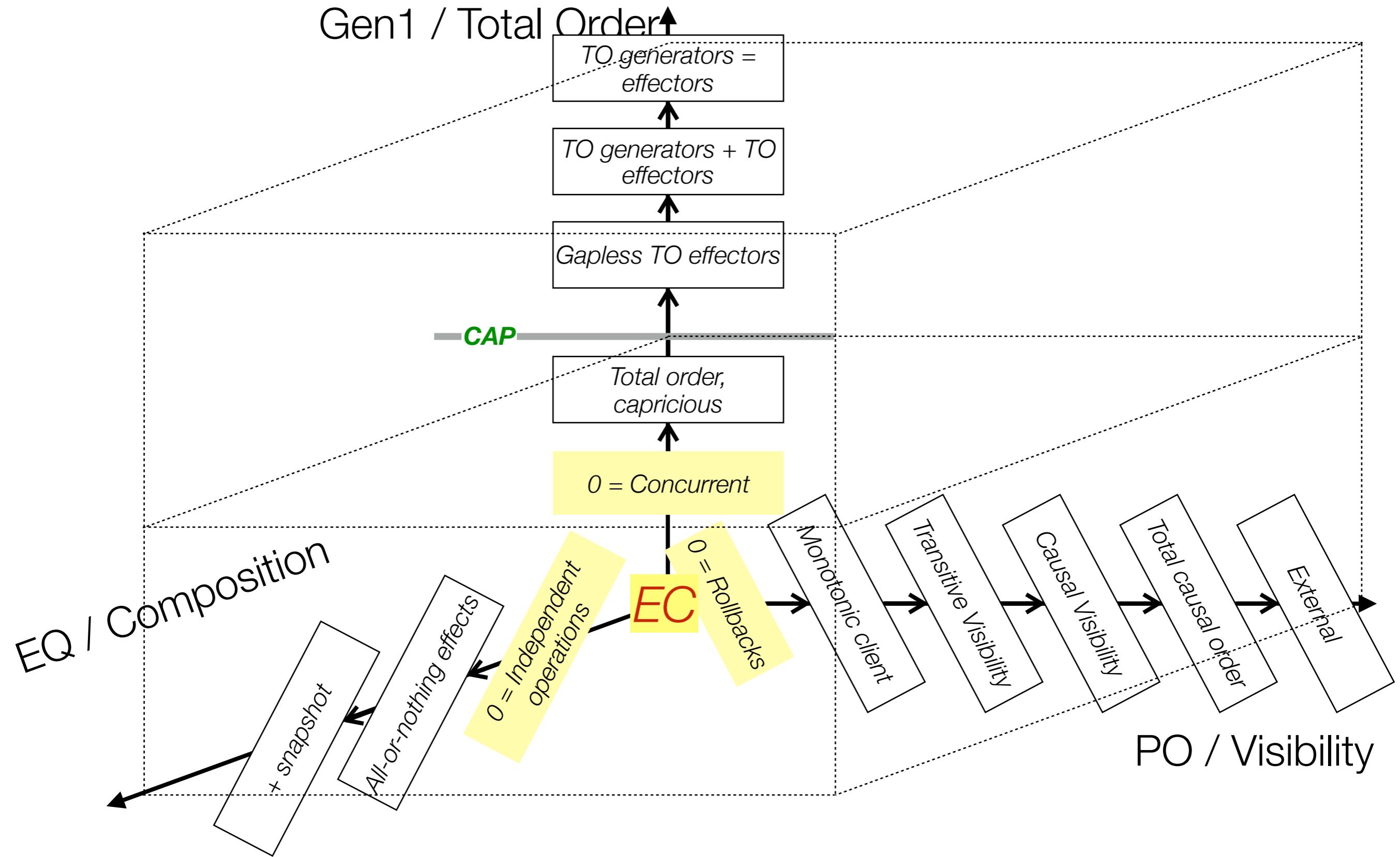
- Assumption: Multiple (n) Objects
- Transactions
 - Write-atomicity: All-or-nothing
 - Read-atomicity: Snapshot

Equality Order Axis

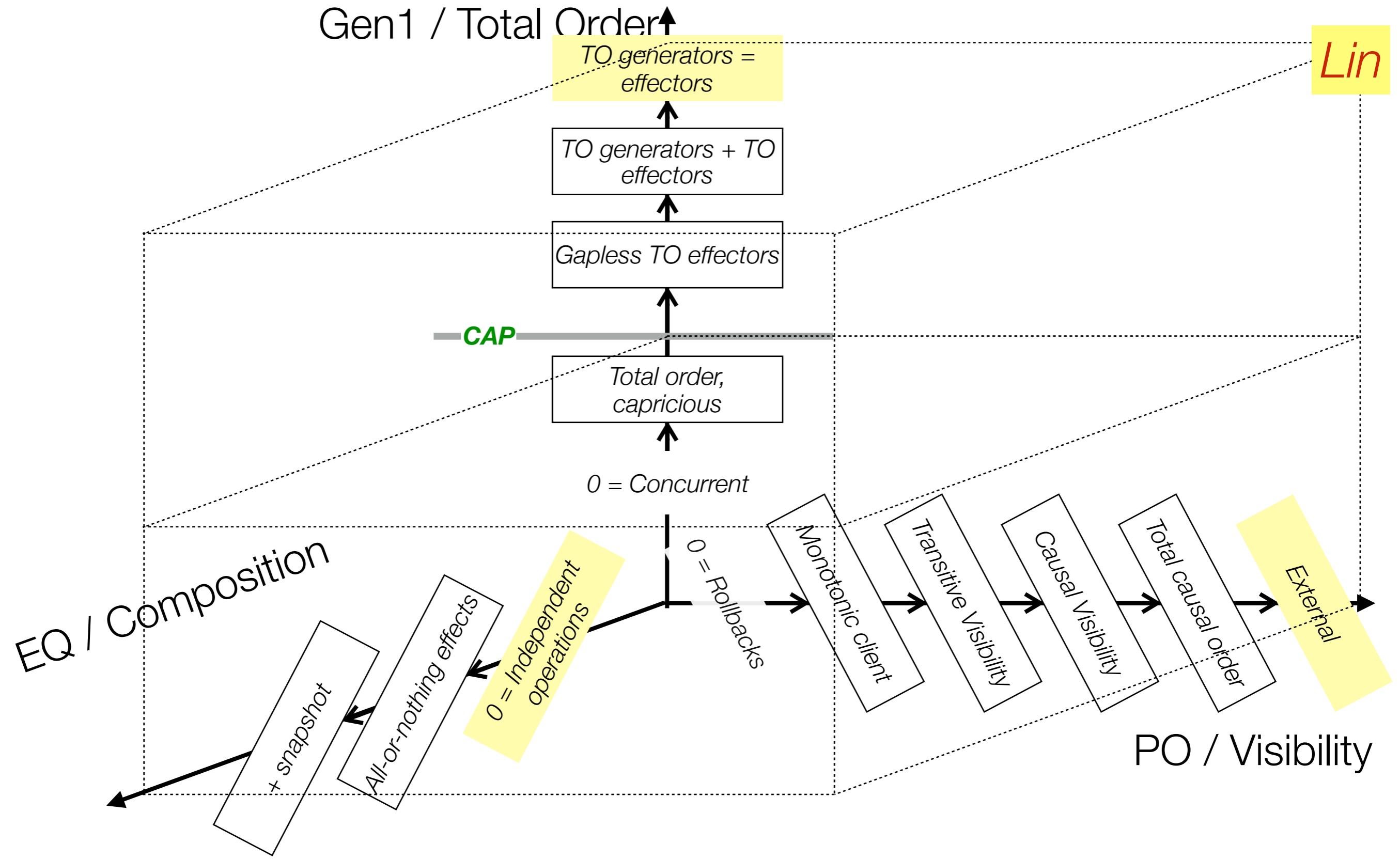
- Assumption: Multiple (n) Objects
- Transactions
 - Write-atomicity: All-or-nothing
 - Read-atomicity: Snapshot
 - Consistent Snapshot



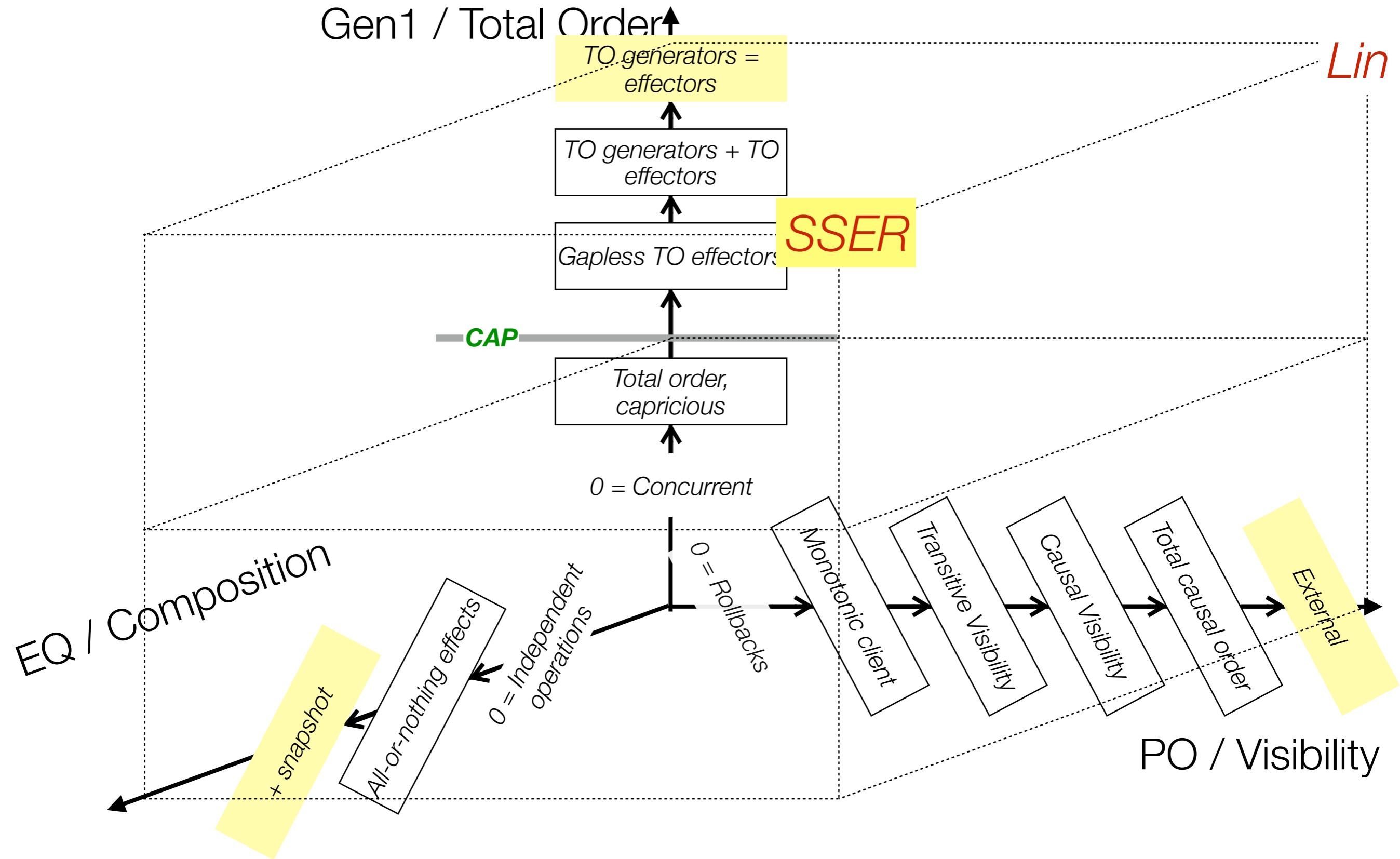
Three dimensions



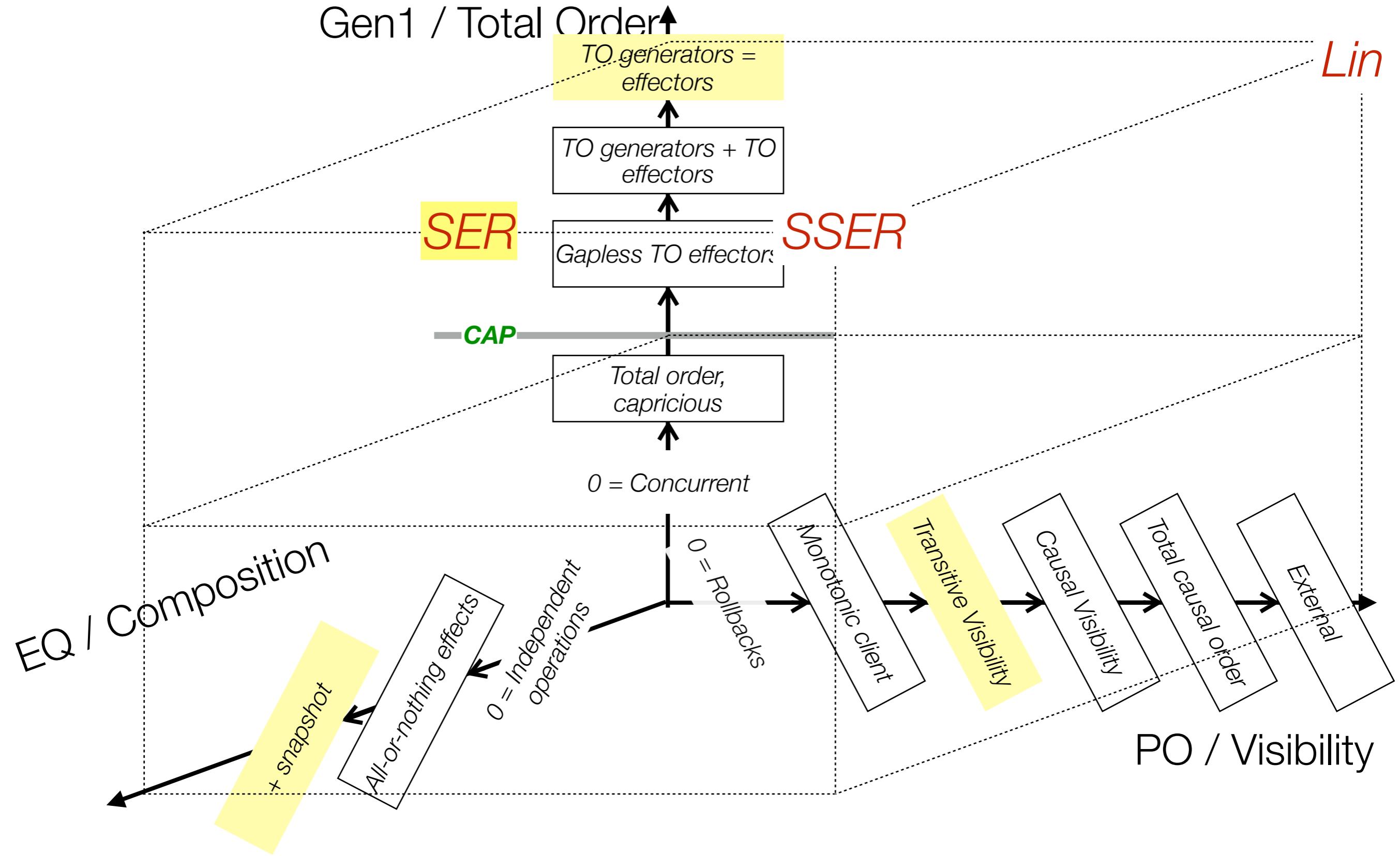
Three dimensions



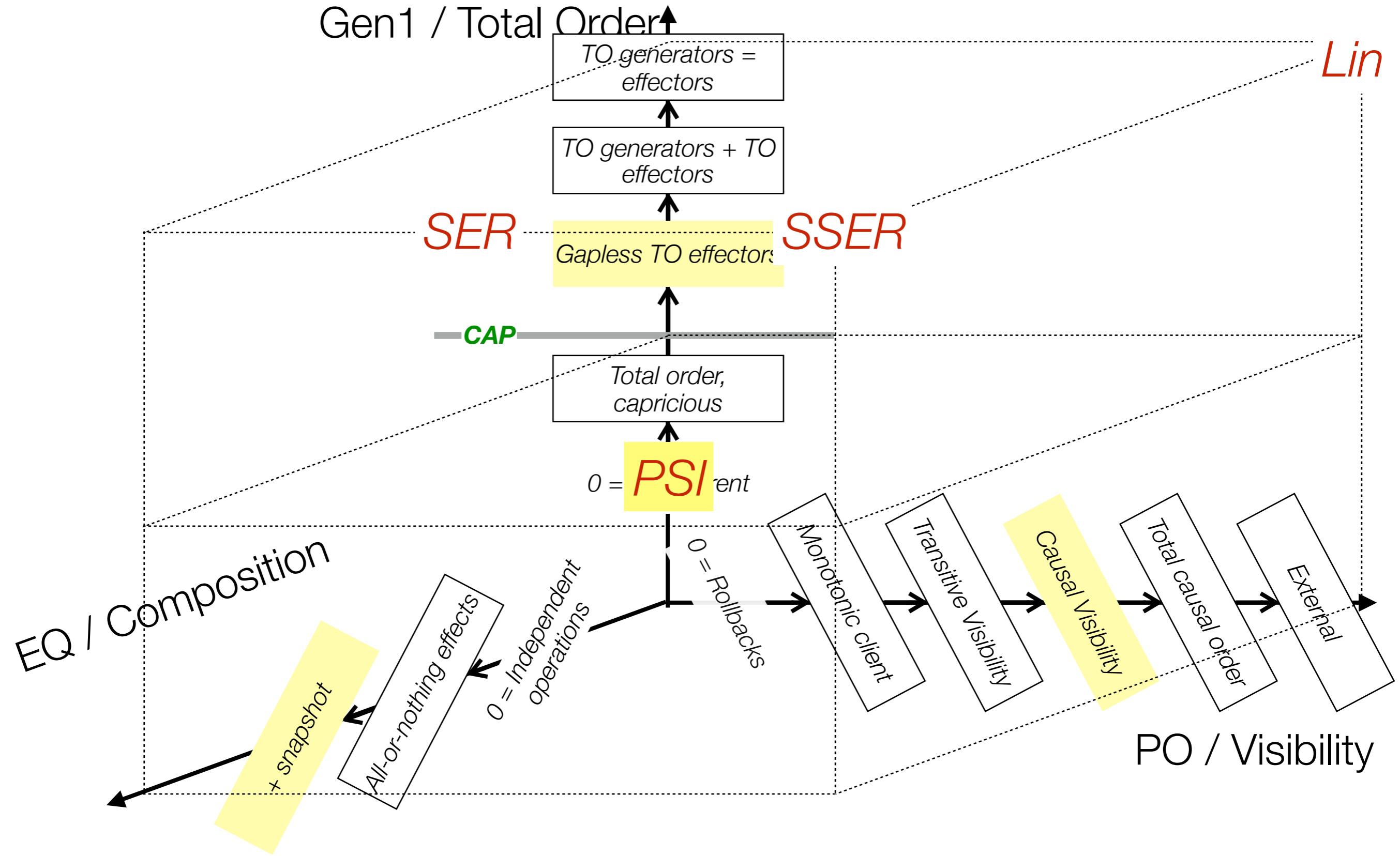
Three dimensions



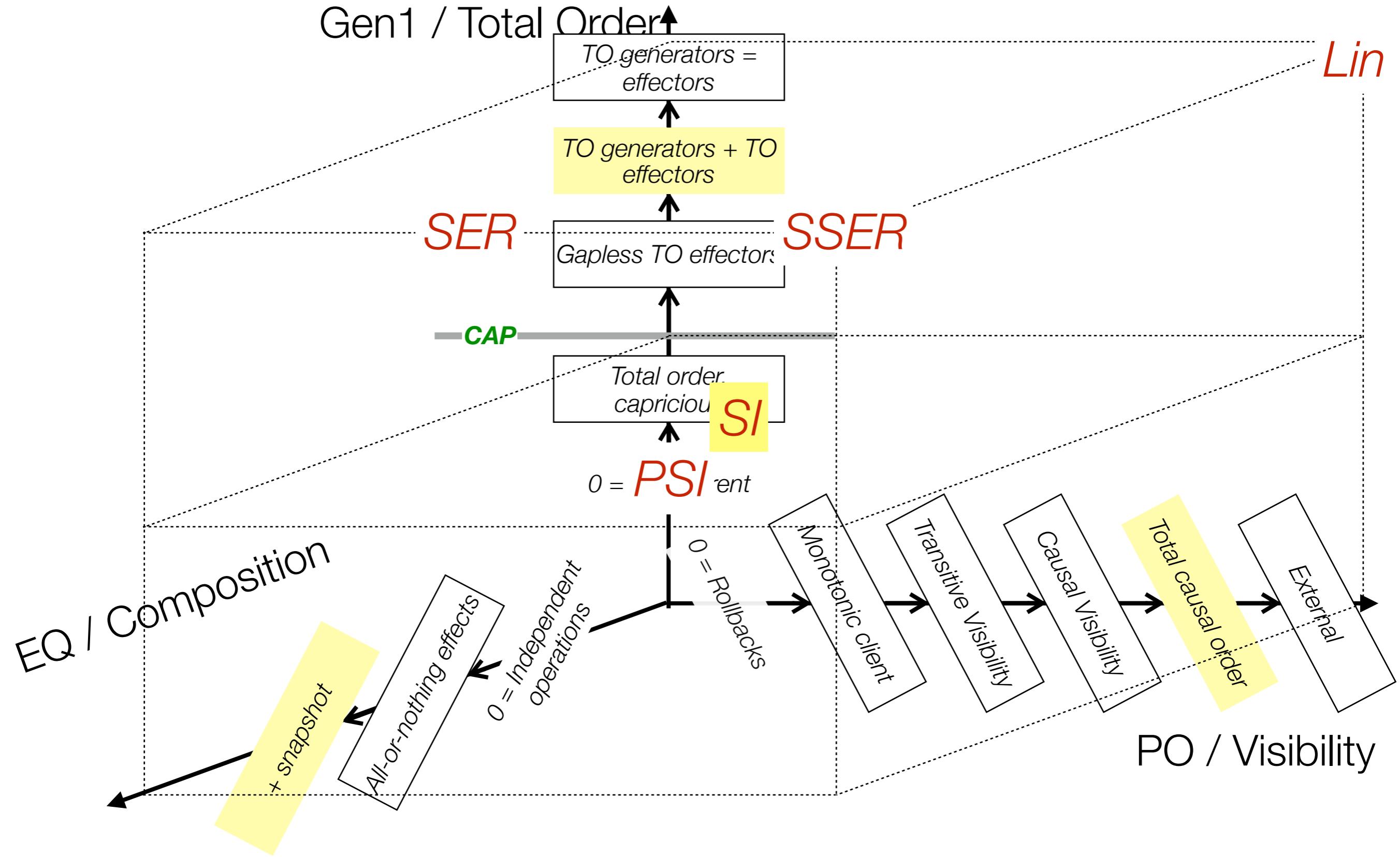
Three dimensions



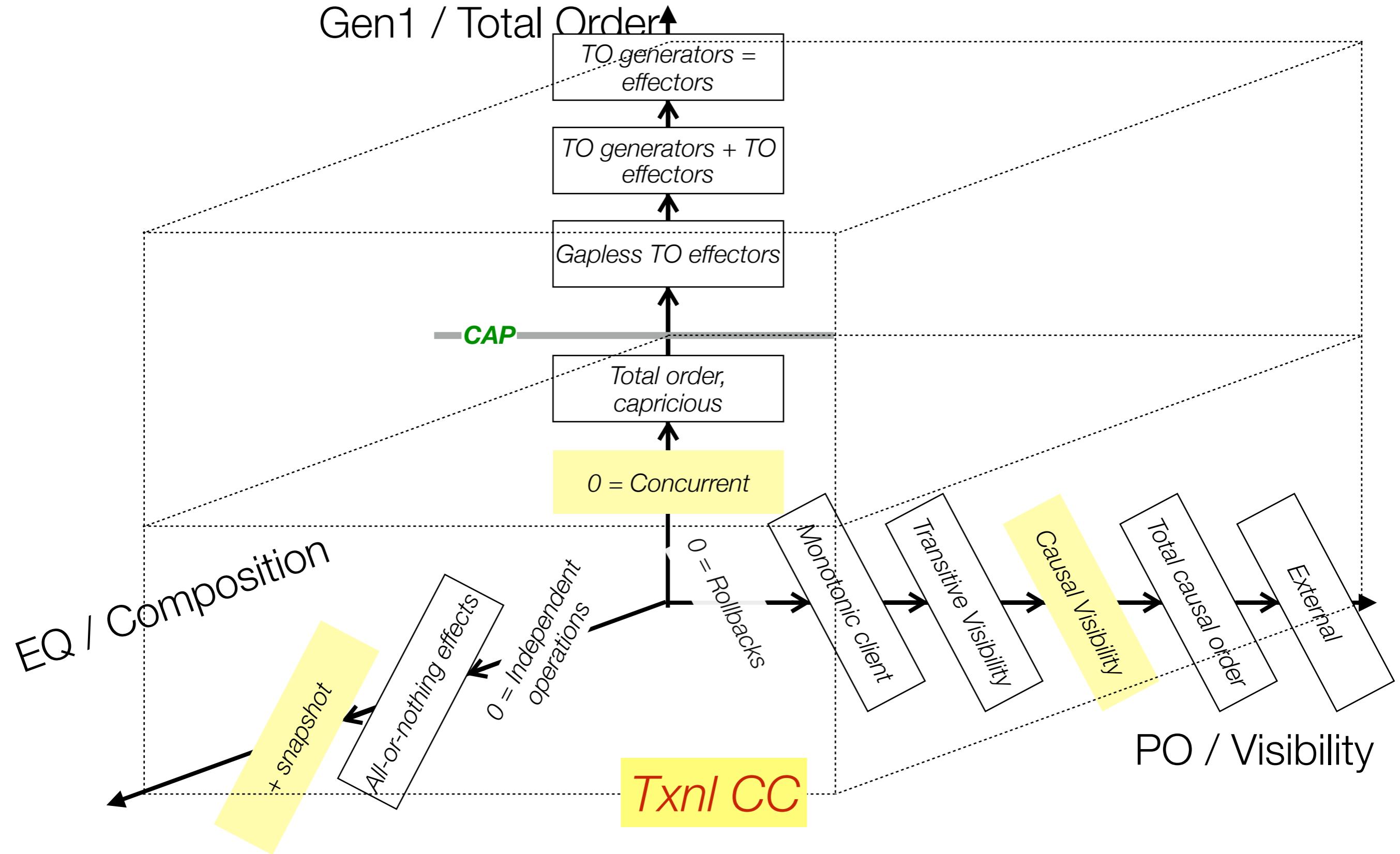
Three dimensions



Three dimensions



Three dimensions



End of day 4