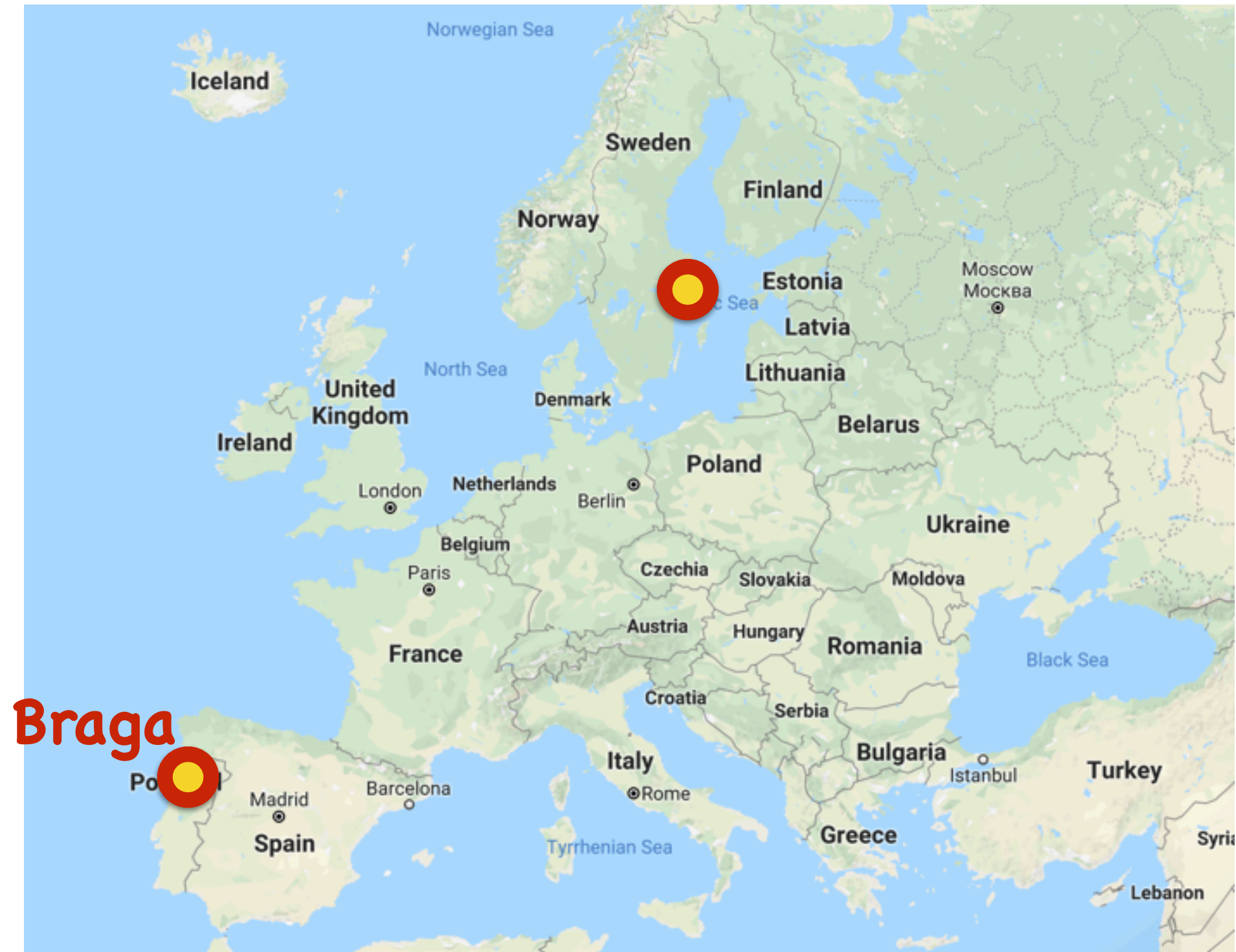




Edge Compute at Hyperscale

Ali Shoker —INESC TEC & Minho univ., Portugal
Stockholm, 2018



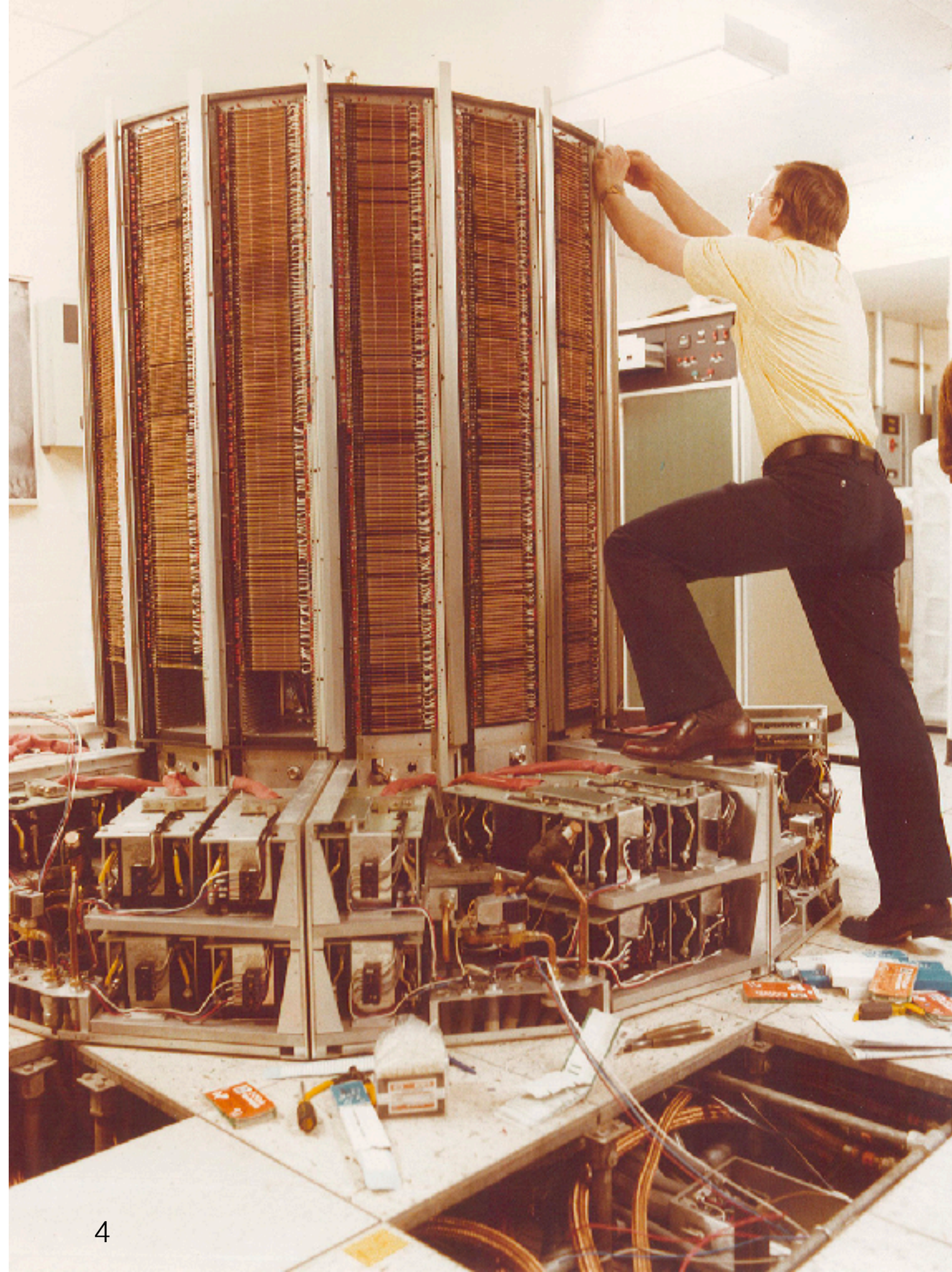
To scale is to scale out

Limit is the sky, but is it

- Cheap?
- Accessible?
- Extendable?
- Maintainable?
- Safe?



Cray-1 Supercomputer (at EPFL), 1975



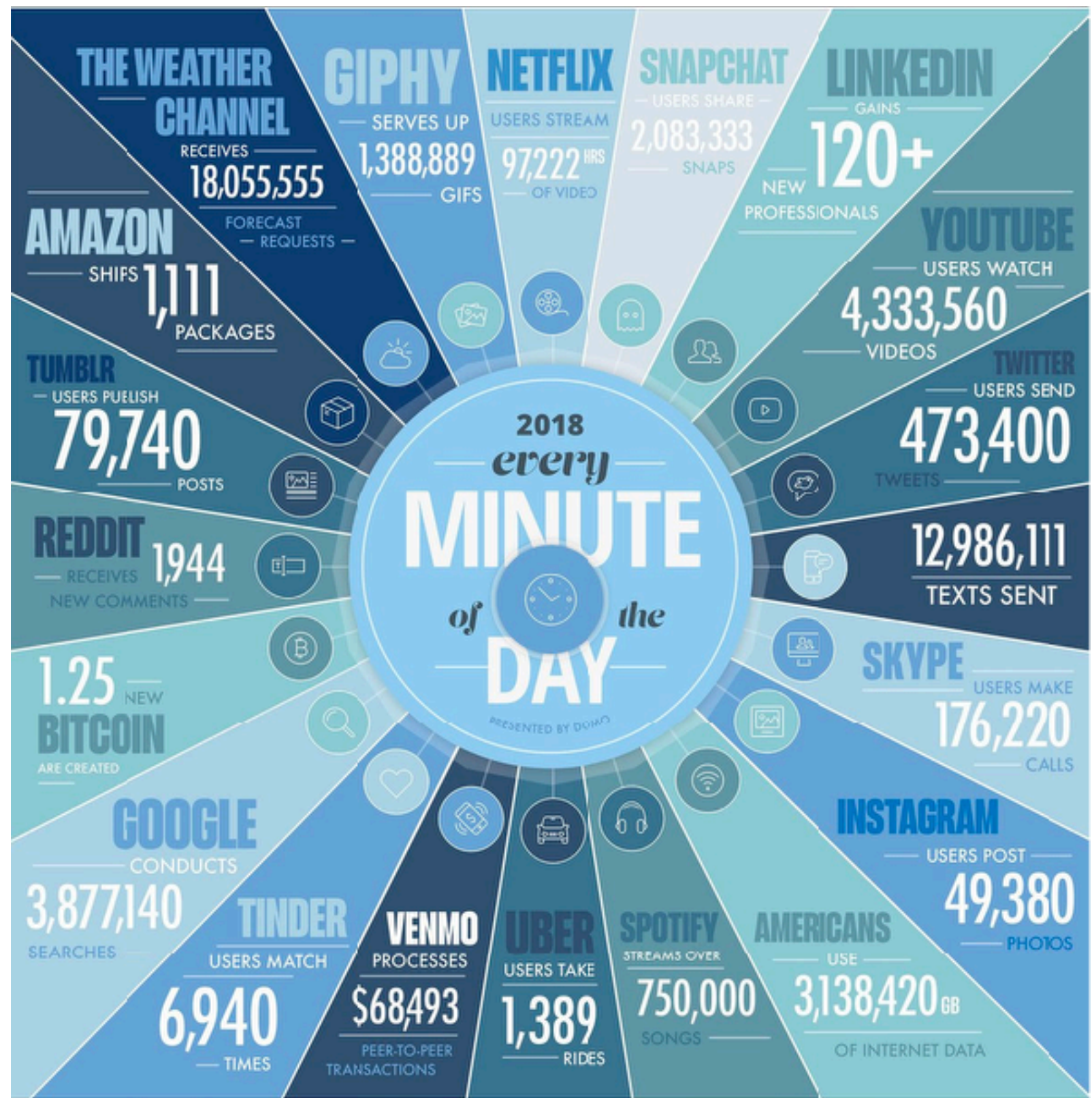
Grid/Cluster computing 90s



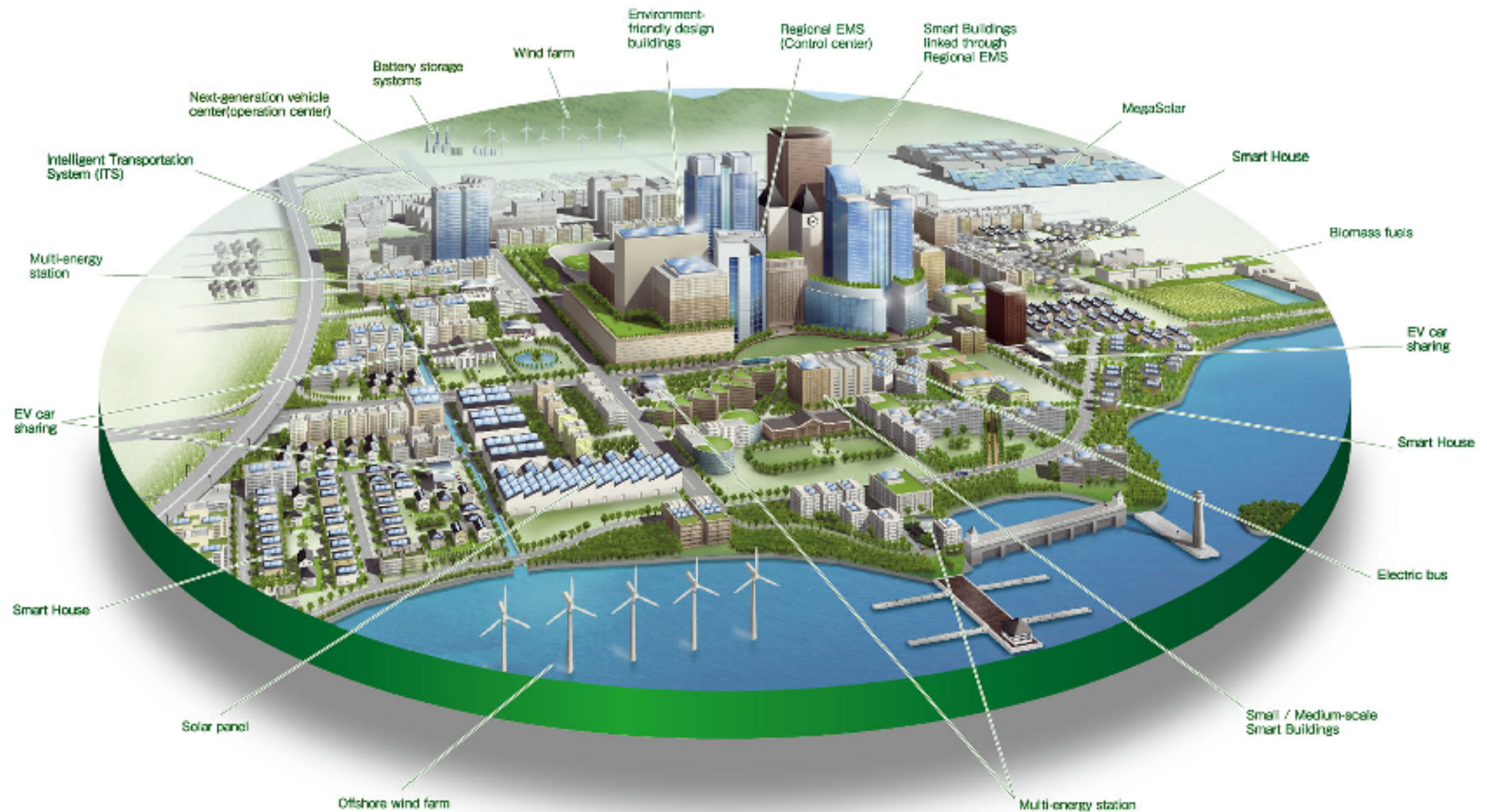
Cloud computing 2000



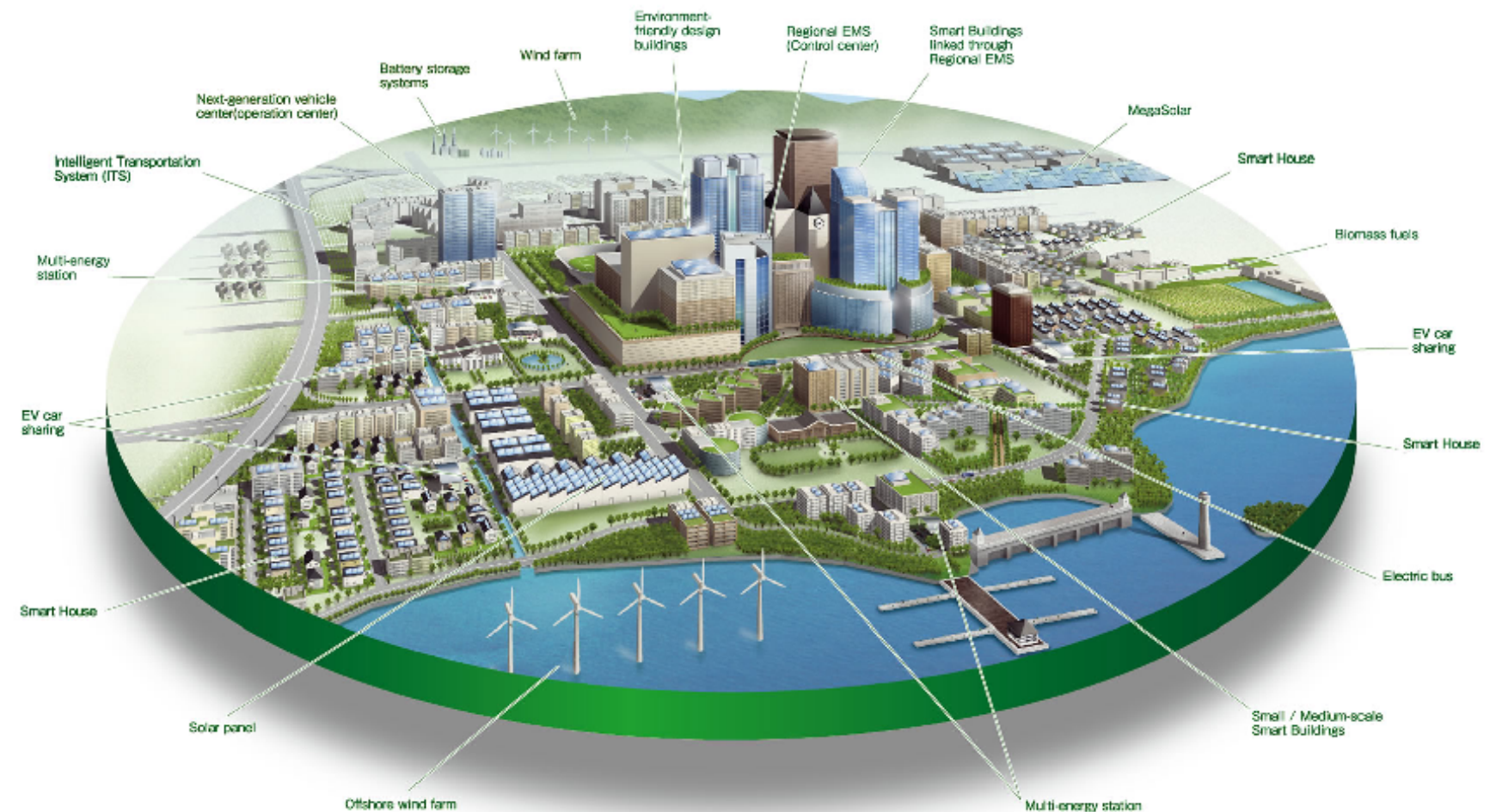
Cloud Overflow



Edge Computing—distributed clouds?



Edge Computing—distributed clouds?



Advantages

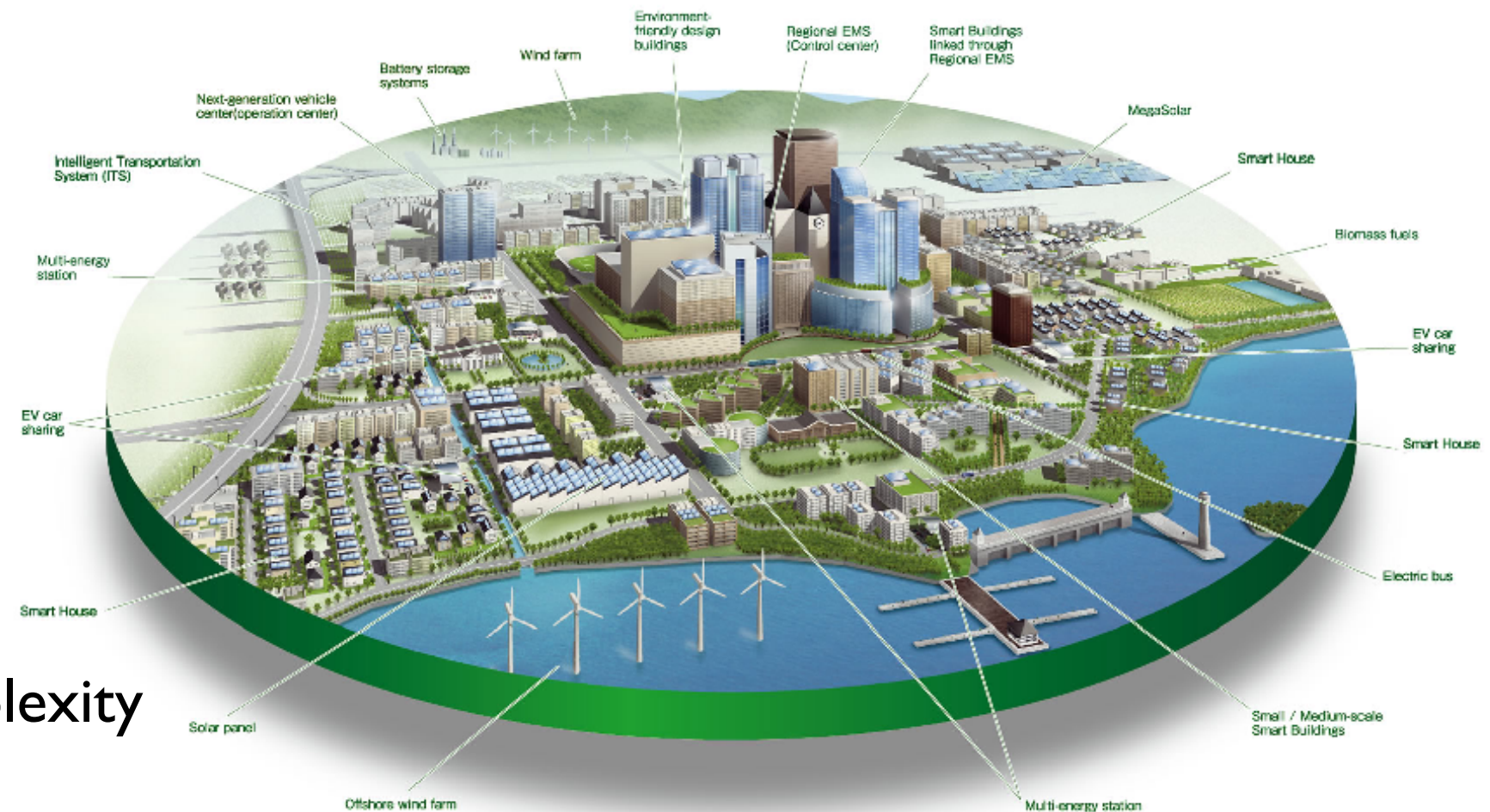
Avoid bottlenecks

Low latency

Save bandwidth

Privacy

Edge Computing—distributed clouds?



Disadvantages

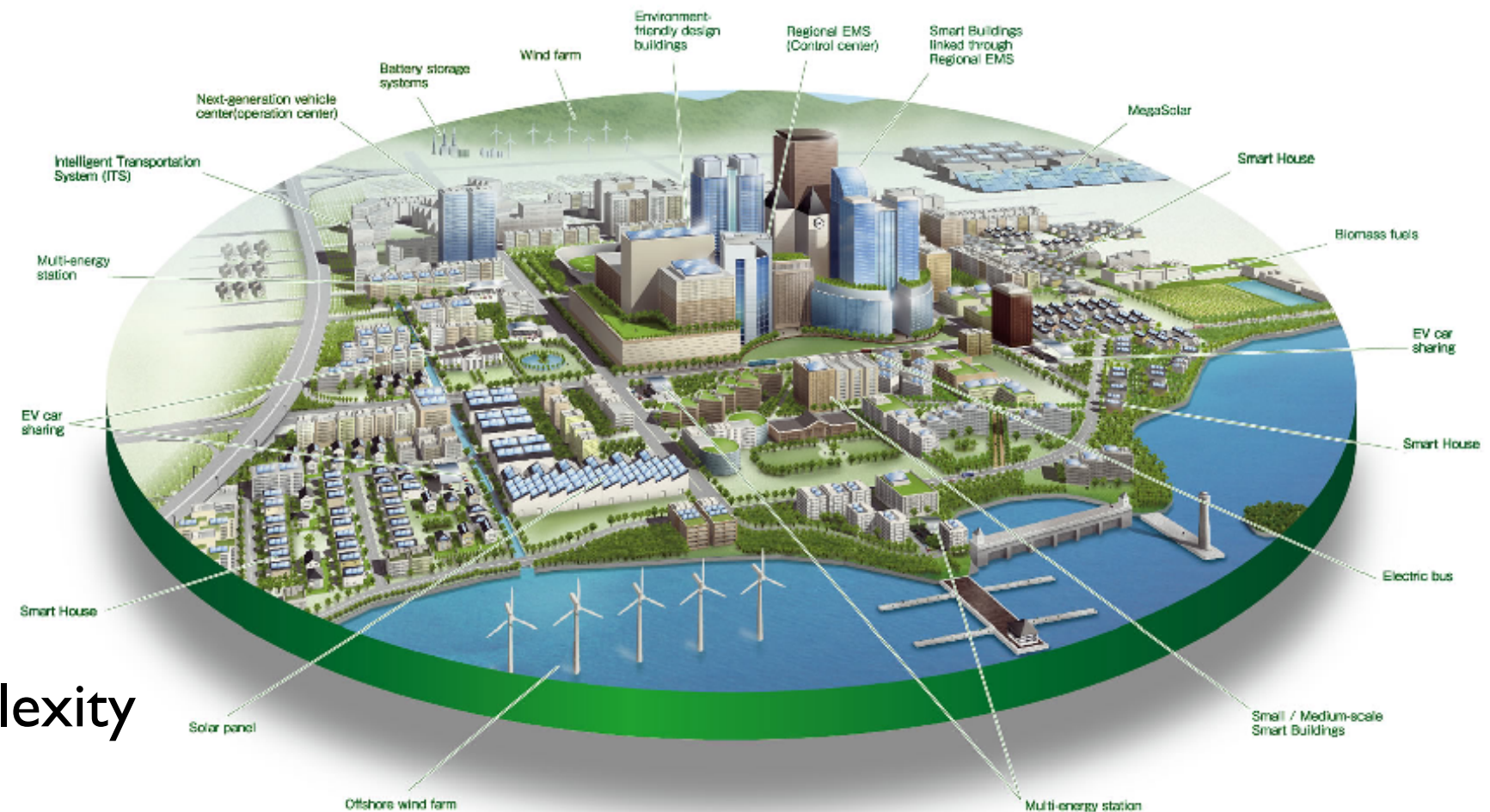
Distributed data complexity

Dynamics

Fragile networks

Heterogeneity...

Edge Computing—distributed clouds?



Disadvantages

Distributed data complexity

Dynamically

Fragile networks

Heterogeneity...

How to build applications?

We focus on data for Edge Computing

Roadmap

How to scale out data

Conflict-free replicated DataTypes (CRDTs)

Hyperscale CRDTs

Outreach and conclusions

How to scale out data?

Partial replication

Full replication

Almost infinite scalability

partition data in **independent entities** with separate scope of serialisability

Life beyond Distributed Transactions: an Apostate's Opinion

Position Paper

Pat Helland

Amazon.Com
705 Fifth Ave South
Seattle, WA 98104
USA

PHelland@Amazon.com

The positions expressed in this paper are personal opinions and do not in any way reflect the positions of my employer Amazon.com.

ABSTRACT

Many decades of work have been invested in the area of distributed transactions including protocols such as 2PC, Paxos, and various approaches to quorum. These protocols provide the application programmer a façade of global

Instead, applications are built using different techniques which do not provide the same transactional guarantees but still meet the needs of their businesses.

This paper explores and names some of the practical approaches used in the implementations of large-scale mission-critical applications in a world which rejects distributed transactions. We discuss the management of fine-grained pieces of application data which may be repartitioned over time as the application grows. We also discuss

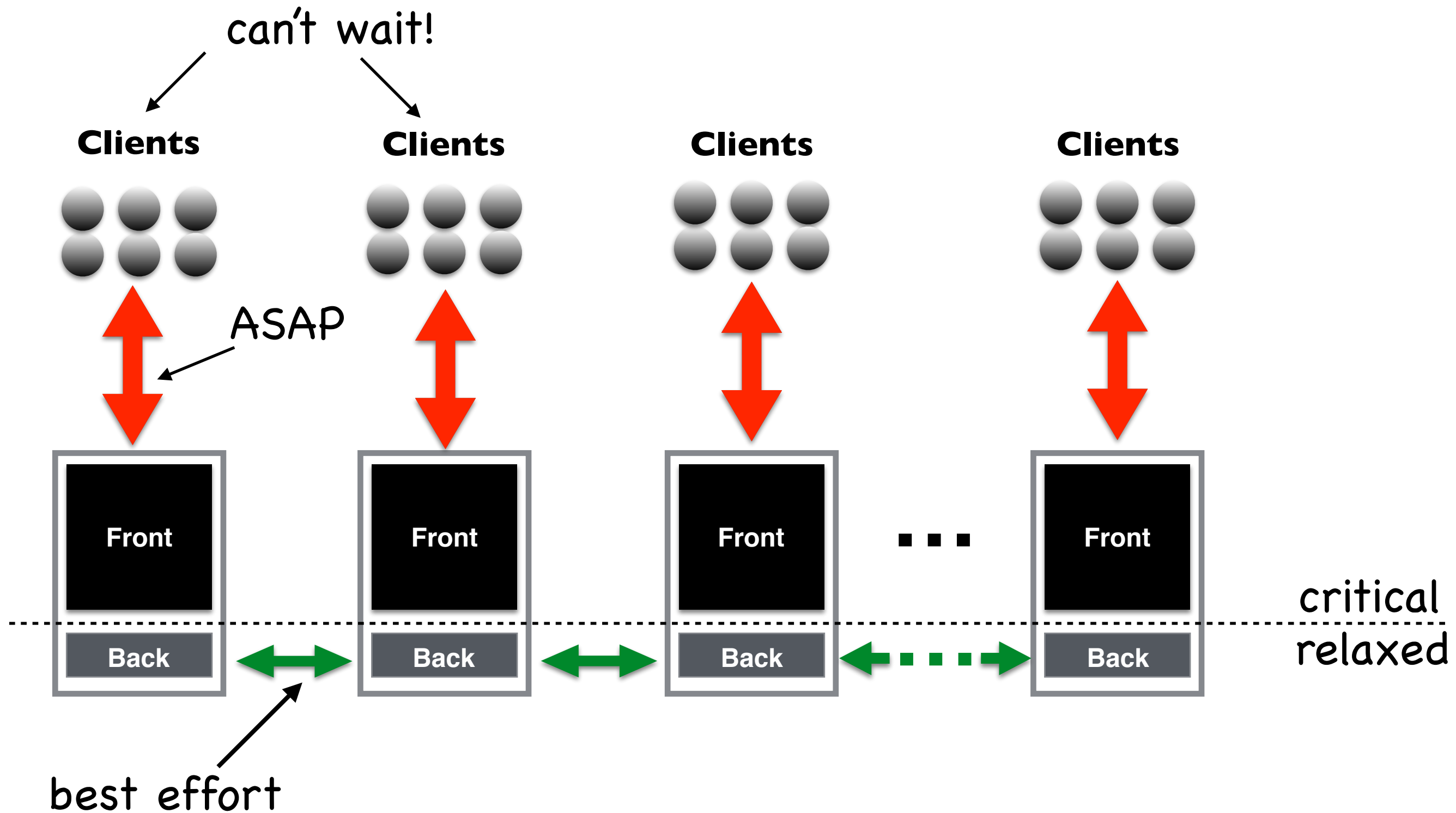
Full replication

CAP theorem: impossible to guarantee all the three:
Consistency, **A**vailability and **P**artition Tolerance?

Maybe use **relaxed consistency**?

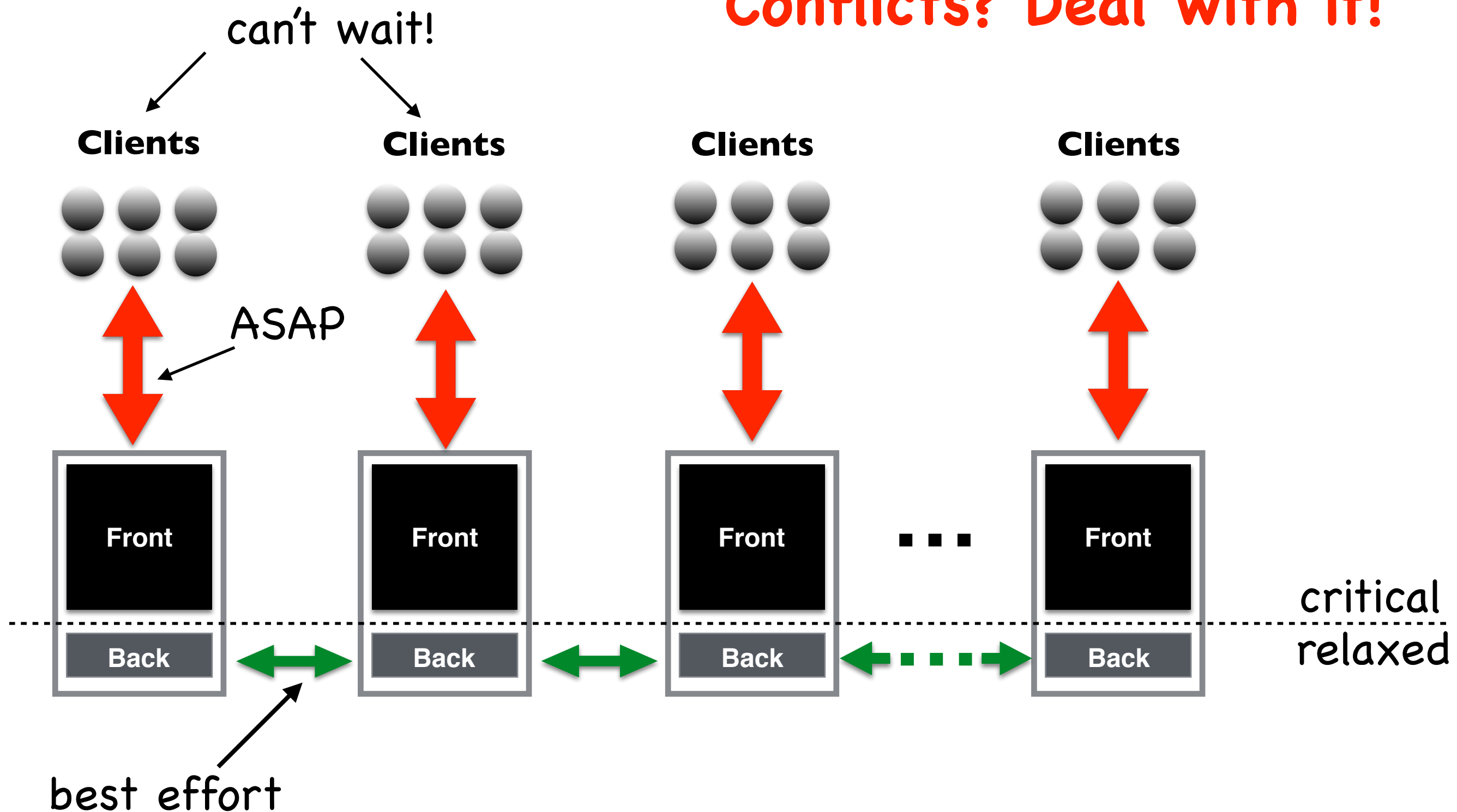


Relaxed consistency



Relaxed consistency

Conflicts? Deal with it!



Conflict-free Replicated DataTypes (CRDT)

What are they?

Data Types: counters, sets, maps, registers, sequences, ...

Replicated: fully or partially (relaxed consistency)

Conflict-free: mathematically proven to eventually converge

Properties

low sync, commutative operations, internal meta-data

Variants

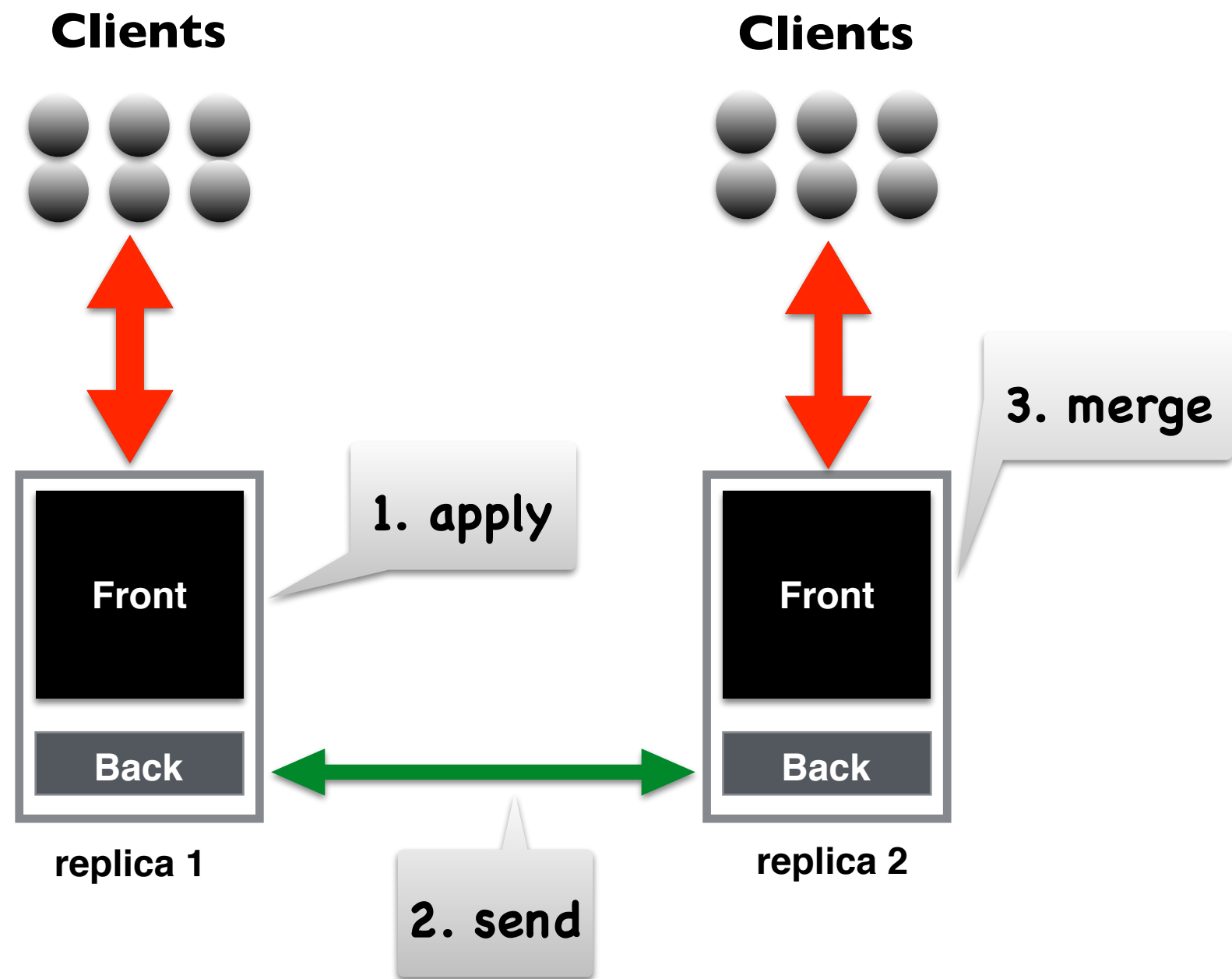
State-based: disseminate a state

Operation-based: disseminate an operation

State-based CRDT

Workflow

1. apply operation
2. send state (in background)
3. merge state remotely



State-based CRDT

Workflow

1. apply operation
2. send state (in background)
3. merge state remotely

inputs:

| $n_i \in \mathcal{P}(\mathbb{I})$, set of neighbors

durable state:

| $X_i := \perp \in S$, CRDT state

on receive _{j,i} (Y)

| $X'_i = X_i \sqcup Y$

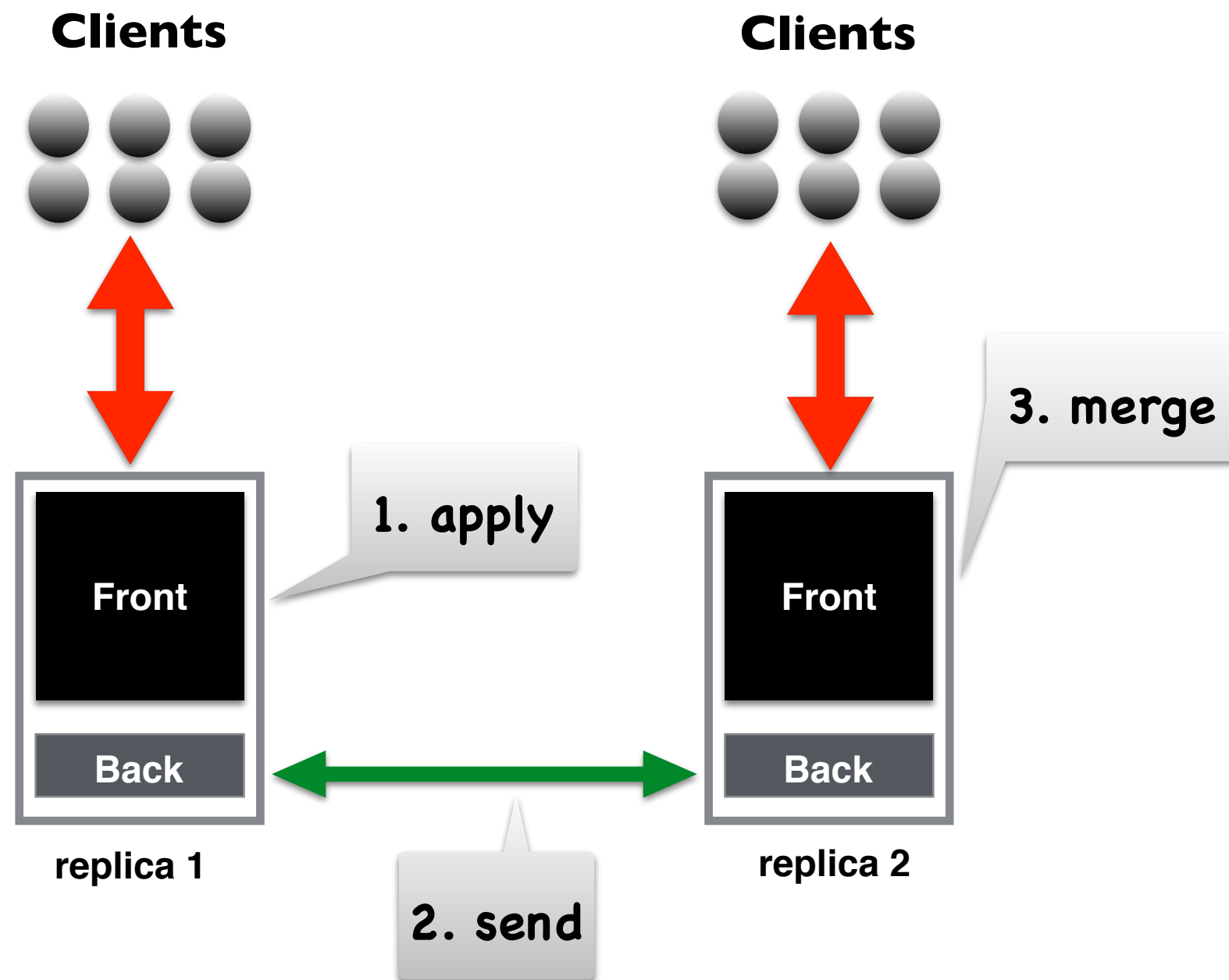
on operation _{i} (m)

| $X'_i = m(X_i)$

periodically // ship state

| $j = \text{random}(n_i)$

| send _{i,j} (X_i)



State-based CRDT

Σ	: State defined as <i>join semi-lattice</i> , σ_i is an instance
<i>Mutators</i>	: mutating operations that inflate the state.
$s \sqcup s'$: LUB to merge states s and s'
$\text{val}_i(q, \sigma_i)$: Read-only evaluation of query q on a state

State-based CRDT

Σ	: State defined as <i>join semi-lattice</i> , σ_i is an instance
<i>Mutators</i>	: mutating operations that inflate the state.
$s \sqcup s'$: LUB to merge states s and s'
$\text{val}_i(q, \sigma_i)$: Read-only evaluation of query q on a state

- Marc Shapiro, Nuno Preguiça, Carlos Baquero, Marek Zawirski. A comprehensive study of Convergent and Commutative Replicated Data Types. [Research Report] RR-7506, Inria – Centre ParisRocquencourt; INRIA. 2011, pp.50.
- Paulo Sérgio Almeida, Ali Shoker, and Carlos Baquero. "Delta state replicated data types". JPDC Elsevier journal, 2018.
- Carlos Baquero, Paulo Sérgio Almeida, and Ali Shoker. "Pure Operation-Based Replicated Data Types". arXiv CoRR . October, 2017.

Grow-only Counter (GCounter) CRDT

A shared counter among nodes: n_1 , n_2 , and n_3

On n_2 :

State: a vector $v=(2,0,4)$

Increment: only my index 2: $v=(2, \textcolor{red}{1}, 4)$

Read: the sum $\Rightarrow 2 + 1 + 4 = \textcolor{red}{7}$

Merge: **(2,0,6)** received from n_3 :

Pair-wise-max $\{(2, \textcolor{green}{1}, \textcolor{blue}{4}), (2, \textcolor{green}{0}, \textcolor{blue}{6})\} \rightarrow \textcolor{red}{(2, 1, 6)}$

$$\begin{aligned}\Sigma &= I \hookrightarrow \mathbb{N} \\ \sigma_i^0 &= \{\} \\ \text{inc}_i(m) &= m\{i \mapsto m(i) + 1\} \\ \text{val}_i(m) &= \sum_{r \in \text{dom}(m)} m(r) \\ m \sqcup m' &= \mathbf{max}(m, m')\end{aligned}$$

So far so good, but how scalable is it?

So far so good, how scalable is it?

State grows linear with number of replicas

$$\begin{aligned}\Sigma &= I \hookrightarrow \mathbb{N} \\ \sigma_i^0 &= \{\} \\ \text{inc}_i(m) &= m\{i \mapsto m(i) + 1\} \\ \text{val}_i(m) &= \sum_{r \in \mathbf{dom}(m)} m(r) \\ m \sqcup m' &= \mathbf{max}(m, m')\end{aligned}$$

Almost infinite scalability

partition data in **independent entities** with separate scope of serialisability

Life beyond Distributed Transactions: an Apostate's Opinion

Position Paper

Pat Helland

Amazon.Com
705 Fifth Ave South
Seattle, WA 98104
USA

PHelland@Amazon.com

The positions expressed in this paper are personal opinions and do not in any way reflect the positions of my employer Amazon.com.

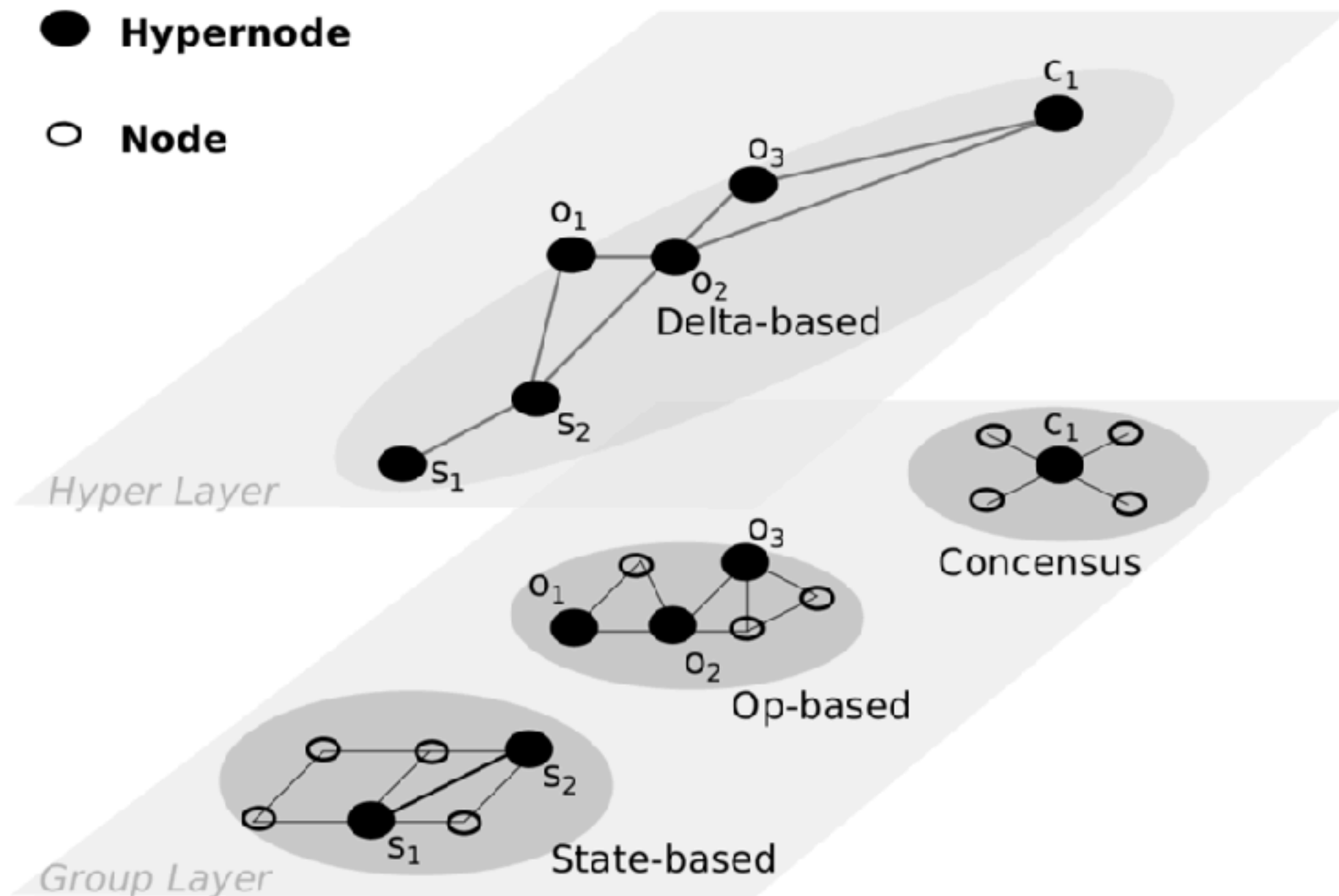
ABSTRACT

Many decades of work have been invested in the area of distributed transactions including protocols such as 2PC, Paxos, and various approaches to quorum. These protocols provide the application programmer a façade of global

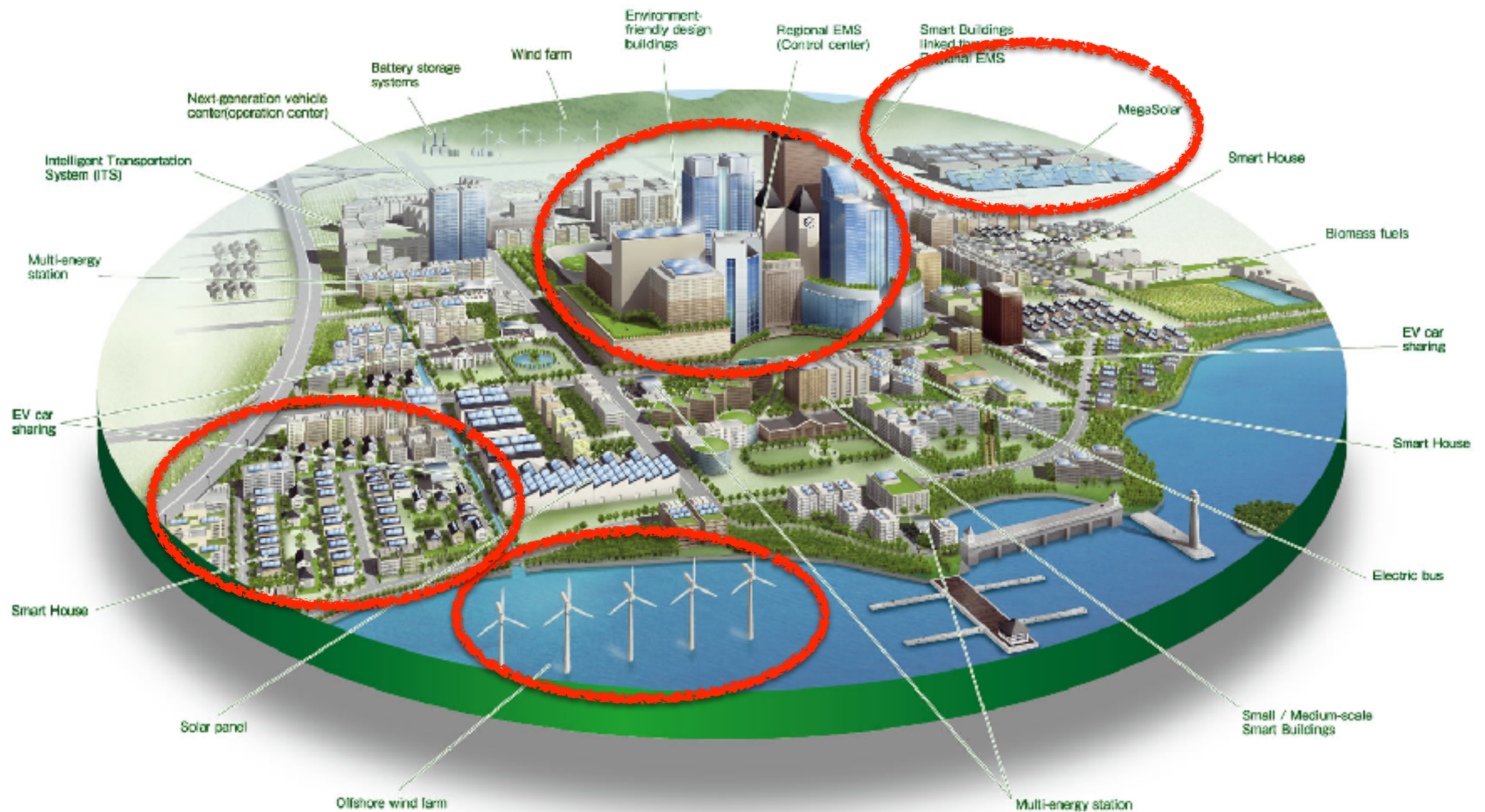
Instead, applications are built using different techniques which do not provide the same transactional guarantees but still meet the needs of their businesses.

This paper explores and names some of the practical approaches used in the implementations of large-scale mission-critical applications in a world which rejects distributed transactions. We discuss the management of fine-grained pieces of application data which may be repartitioned over time as the application grows. We also discuss

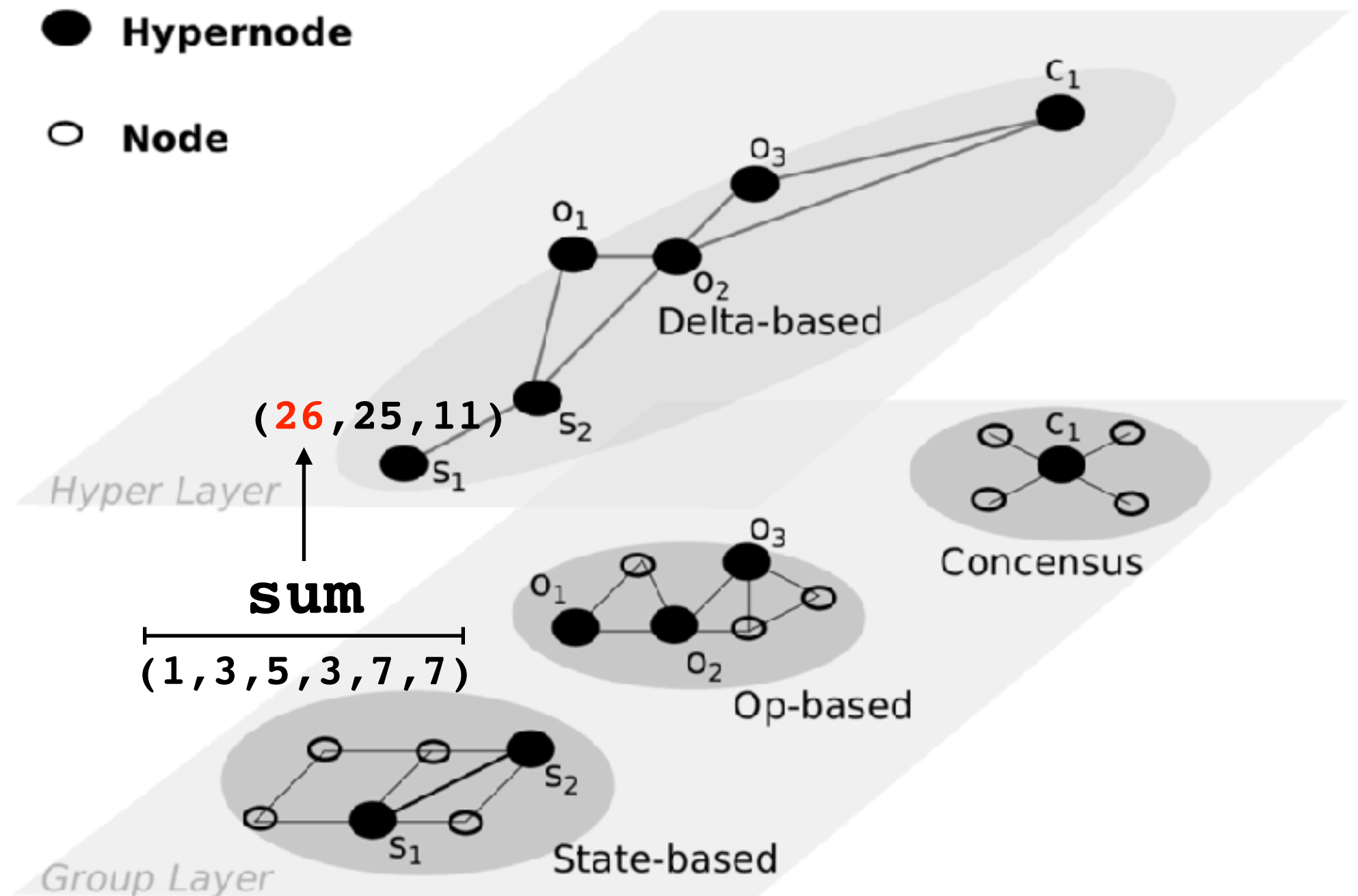
Hyperscale CRDTs



Edge Compute at Hyperscale



Hyperscale GCounter CRDTs



HyperGCounter CRDT

$$\begin{aligned}\text{GCDatum}(\mathbb{I}) &= \mathbb{I} \hookrightarrow \mathbb{N} && // \text{a map } m \\ \perp &= \{\} \\ \text{inc}_i(m) &= m\{i \mapsto m(i) + 1\} \\ \text{datumvalue}_i(m) &= \sum_{j \in \mathbb{I}} m(j) \\ m \sqcup m' &= \text{pair_wise_max}(m, m')\end{aligned}$$

HyperGCounter CRDT

$$\begin{aligned}\text{GCDatum}(\mathbb{I}) &= \mathbb{I} \hookrightarrow \mathbb{N} && // \text{a map } m \\ \perp &= \{\} \\ \text{inc}_i(m) &= m\{i \mapsto m(i) + 1\} \\ \text{datumvalue}_i(m) &= \sum_{j \in \mathbb{I}} m(j) \\ m \sqcup m' &= \text{pair_wise_max}(m, m')\end{aligned}$$

$$\begin{aligned}\text{GCDigest}(\mathbb{G}) &= \mathbb{G} \hookrightarrow \mathbb{N} && // \text{a map } m \\ \perp &= \{\} \\ \text{update}_k(m, V) &= m\{k \mapsto \max(m(k), V)\} \\ \text{digestvalue}_k(m) &= \sum_{j \in \mathbb{G}} m(j) \\ m \sqcup m' &= \text{pair_wise_max}(m, m')\end{aligned}$$

HyperGCounter CRDT

$$\text{GCDatum}(\mathbb{I}) = \mathbb{I} \hookrightarrow \mathbb{N} \quad // \text{a map } m$$

$$\perp = \{\}$$

$$\text{inc}_i(m) = m\{i \mapsto m(i) + 1\}$$

$$\text{datumvalue}_i(m) = \sum_{j \in \mathbb{I}} m(j)$$

$$m \sqcup m' = \text{pair_wise_max}(m, m')$$

$$\text{GCDigest}(\mathbb{G}) = \mathbb{G} \hookrightarrow \mathbb{N} \quad // \text{a map } m$$

$$\perp = \{\}$$

$$\text{update}_k(m, V) = m\{k \mapsto \max(m(k), V)\}$$

$$\text{digestvalue}_k(m) = \sum_{j \in \mathbb{G}} m(j)$$

$$m \sqcup m' = \text{pair_wise_max}(m, m')$$

$$\text{HyperGC}(\mathbb{L}, \mathbb{G}) = \text{GCDatum}(\mathbb{L}) \times \text{GCDigest}(\mathbb{G}) \quad // \text{two maps: } l, g$$

$$\perp = (\perp, \perp)$$

$$\text{inc}_i((l, g)) = (\text{inc}_i(l), \perp) \quad // \text{increment local datum}$$

$$\text{publish}_k((l, g)) = (\perp, \text{update}_k(g, \text{datumvalue}(l))) \quad // \text{only hypervisors}$$

$$\text{hypervalue}_k((l, g)) = \text{digestvalue}_k(g)$$

$$(l, g) \sqcup (l', g') = (l \sqcup l', g \sqcup g')$$

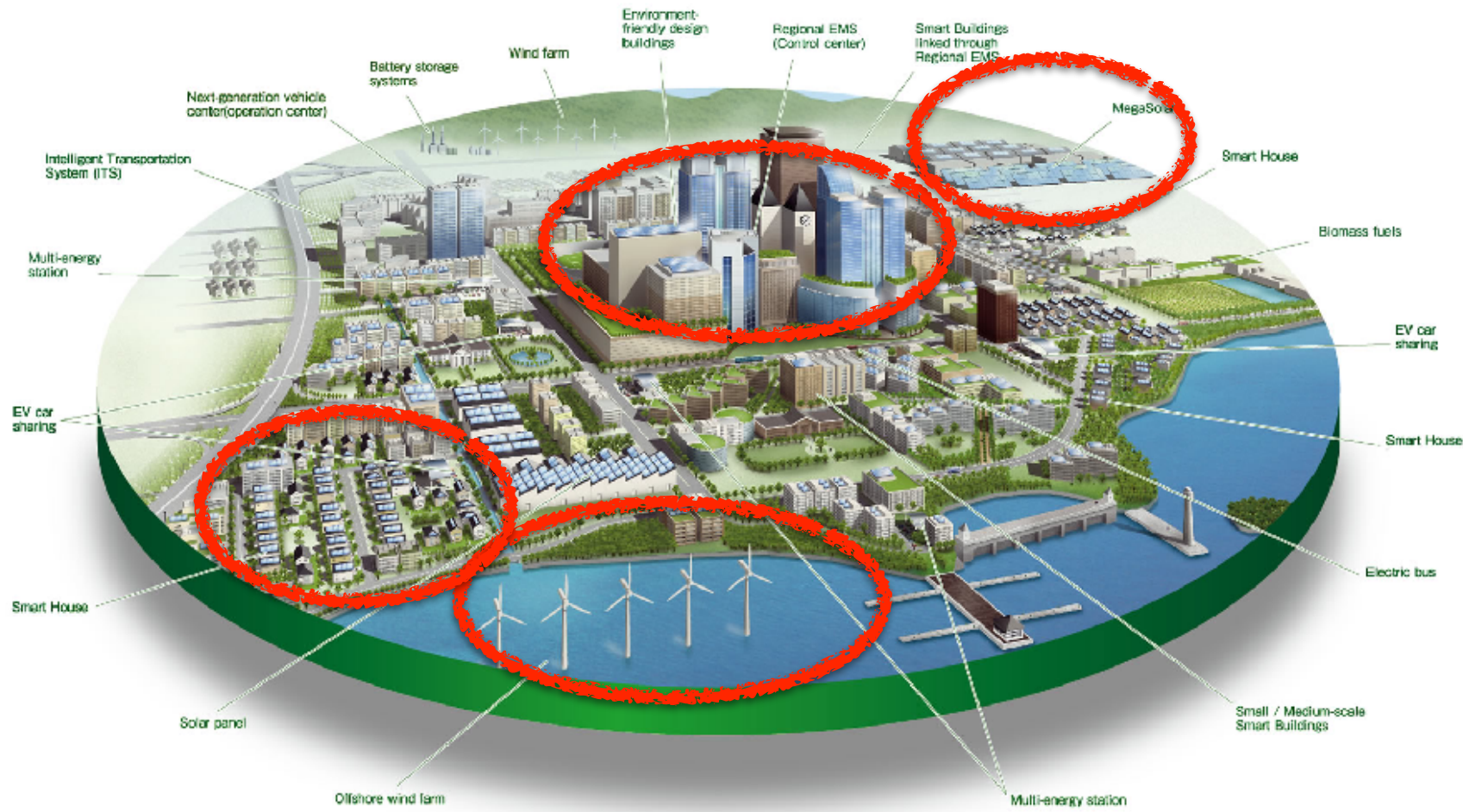
Other Hyper CRDTs

- ☐ GCounter
- ☐ Counter
- ☐ Multi-sets (implemented as maps of counters)
- ☐ Sets (with version-vector metadata)
- ☐ Others

In progress

- ☐ Investigate other CRDTs
- ☐ Allow for remote updates
- ☐ Dynamic membership
- ☐ Empirical evaluations

Concluding remarks



- ❑ To scale is to scale out
- ❑ Scale out CRDTs: $O(N) \Rightarrow O(\log N)$
- ❑ Strength: interoperability, flexibility, storage, bandwidth
- ❑ Weakness: complexity, limitations

Get in touch

Ali Shoker

Web: www.alishoker.com

Twitter: ashokerCS