

CASCADE: FULL DETAILS

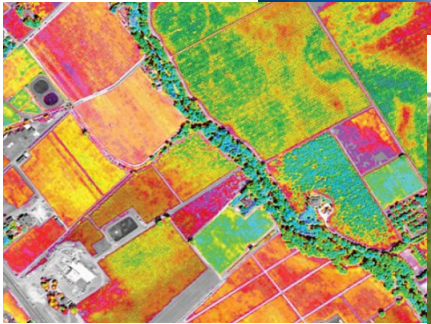
Ken Birman
CS5412 Fall 2022
Lecture 10

THE PROMISE OF EDGE INTELLIGENCE

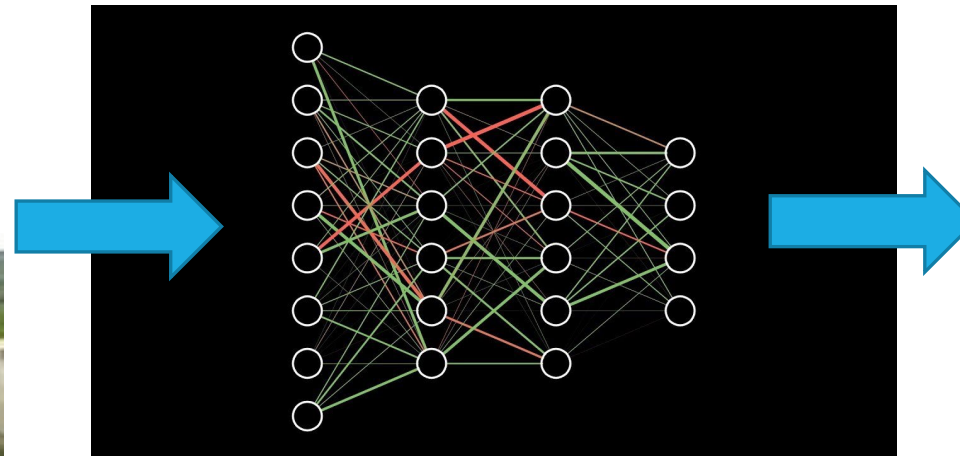
Can “real world” artificial intelligence
be as accurate and as fast as human
intelligence?

THE WORLD IS GENERATING A NEW WAVE OF IOT/ML PIPELINES... THERE ARE MANY USE CASES

Data sources



Cooperative ML
Distributed AI

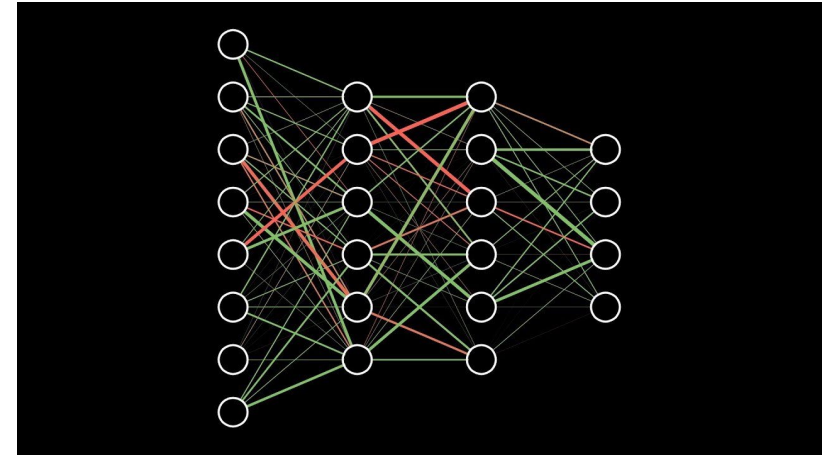


Smart
Queries



**How much should I
budget for raw milk
purchases in March for
my yoghurt factory?"**

COOPERATIVE ML?



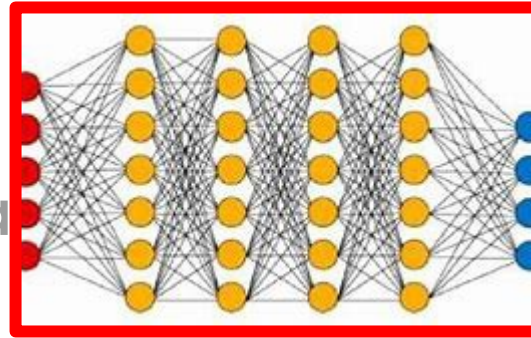
One ML to rule them all? Not likely!

So we need to think about cooperation between MLs in robotics, smart homes, 5G, digital twin scenarios.

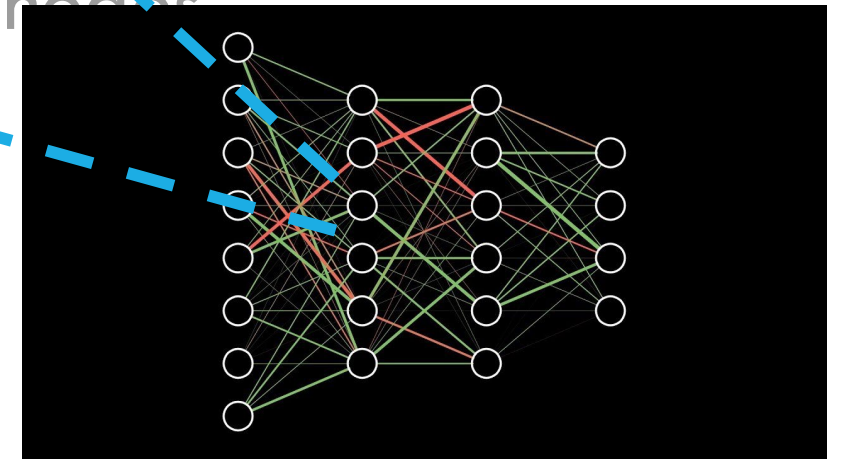
- The application is a graphical collection of AI classifiers / learners
- Nodes represent computational tasks.
- Edges represent data flow between distinct tasks.

A SINGLE λ COULD BE AN ENTIRE DISTRIBUTED AI, ON A COLLECTION OF MACHINES

Single λ may run on a collection of nodes



Today's cloud platforms have limited support for this model, lack the real-time and consistency guarantees needed for IoT.



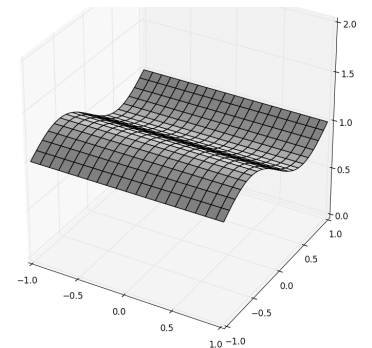
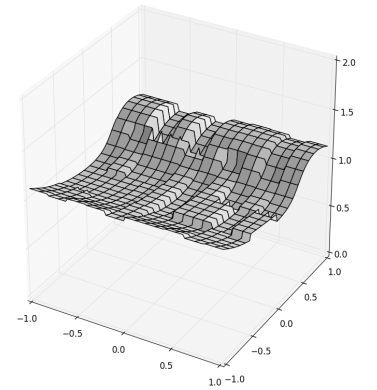
THE NEED IS FOR TIME-FOCUSED EDGE INTELLIGENCE

In today's cloud, AIs
fight platform noise

IoT data arrives as a stream.

Edge intelligence must be *instantly reactive*

- Get power grid data at 250ms intervals
- Format it as 15 x 15 tensor by GPS location with a time axis from 10:00:00 to 10:01:00
- Run ML phase disruption analytic on the tensor



**Cascade: Faster
yet accurate**

HIGH LEVEL CASCADE GOALS



Legacy support: Easy to use with no need to change your code

Much faster than standard platforms: low delay, high bandwidth

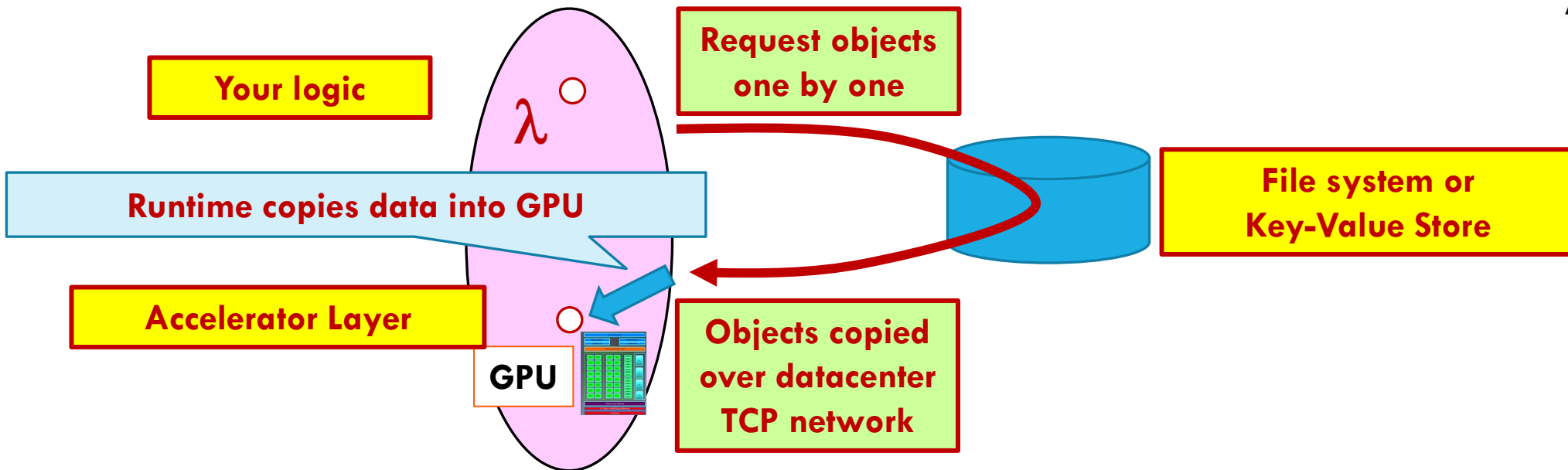
Stronger guarantees: Your ML doesn't fight platform "noise"

BUT HOW WILL WE GUARANTEE CONSISTENCY? HOW CAN IT BE SO FAST?

Let's start by asking how to make it fast

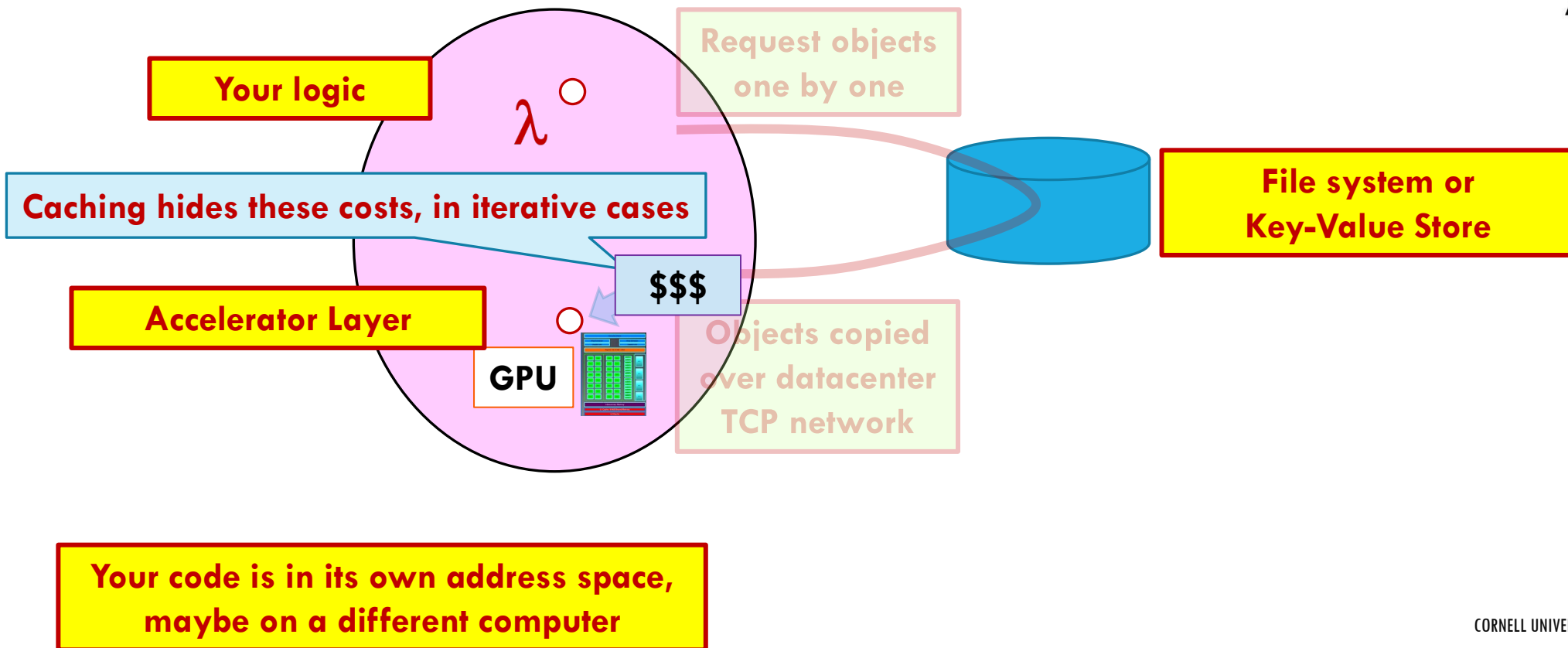
Then we will see that we already know how to guarantee consistency: be smart about time (and time skew), and use Chandy-Lamport consistent cuts when we hit the limits of clock precision/accuracy!

TRADITIONAL APPROACH

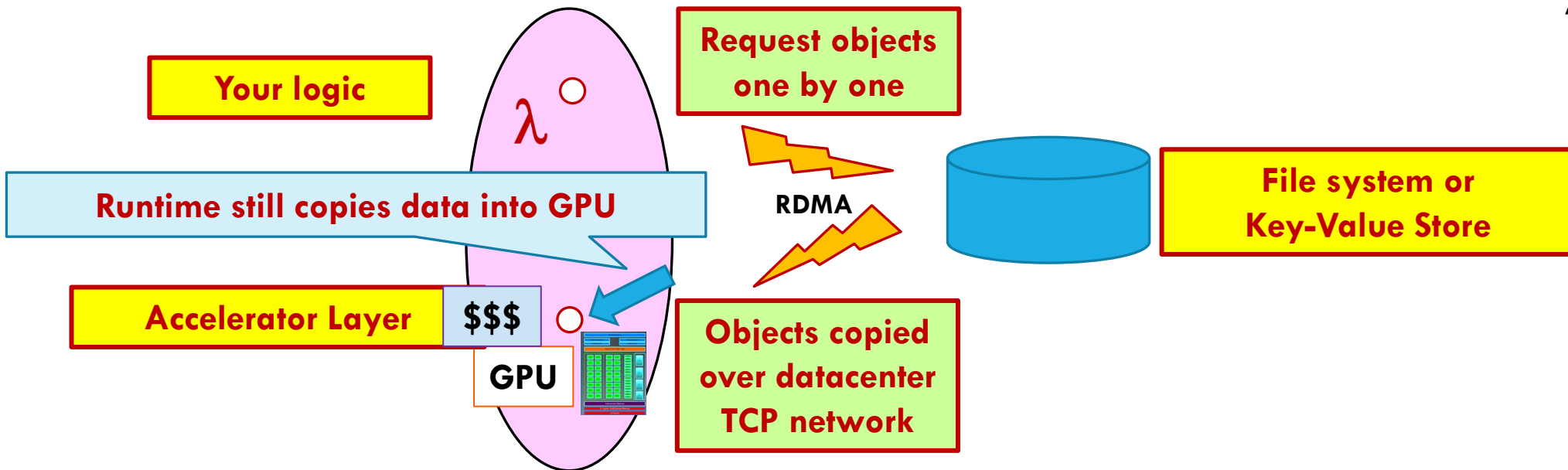


Your code is in its own address space,
maybe on a different computer

TRADITIONAL APPROACH

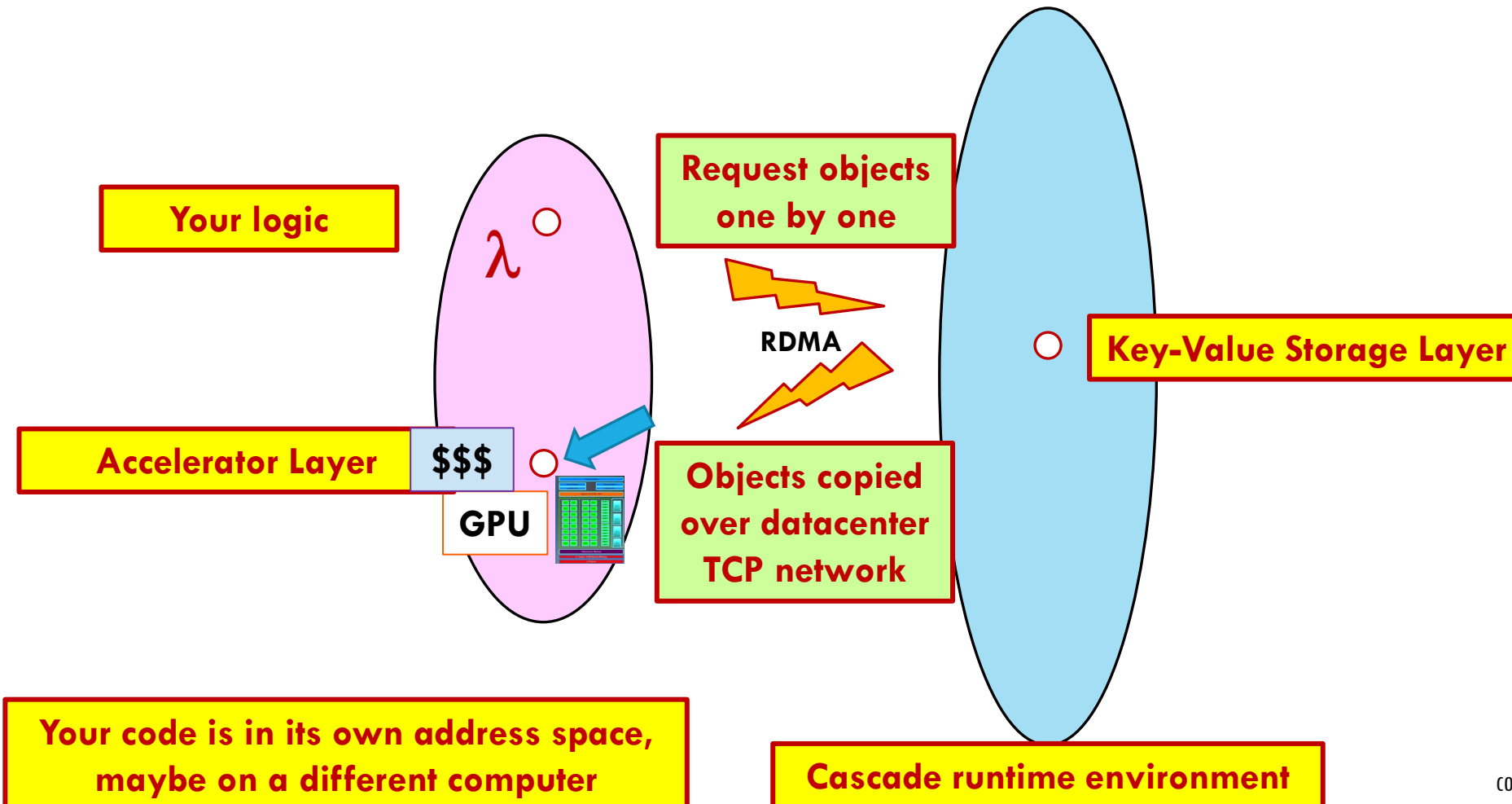
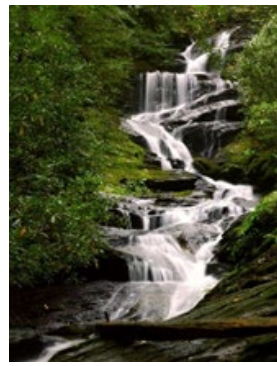


RDMA CAN HELP... A LITTLE

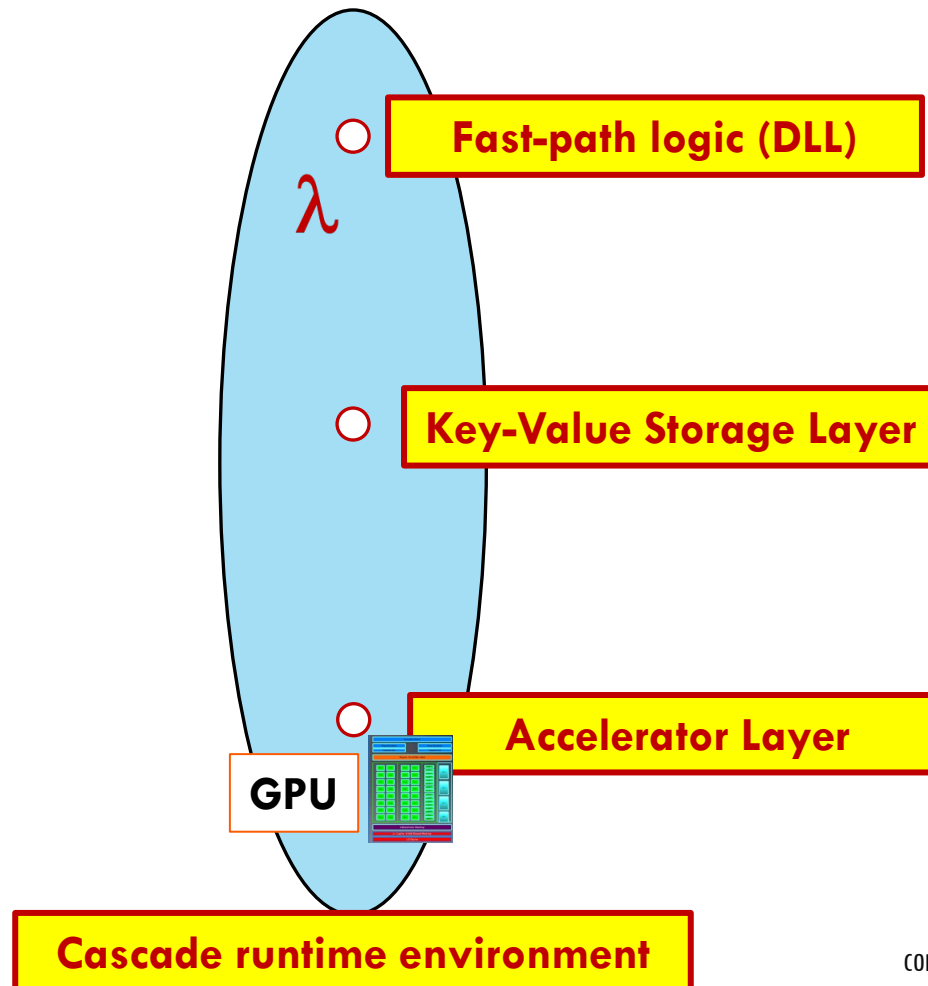
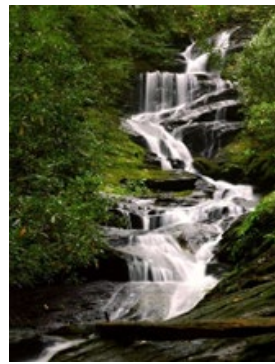


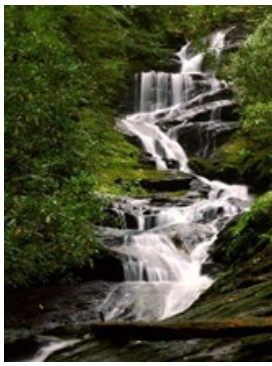
Your code is in its own address space,
maybe on a different computer

CASCADE CAN BE USED THIS WAY...

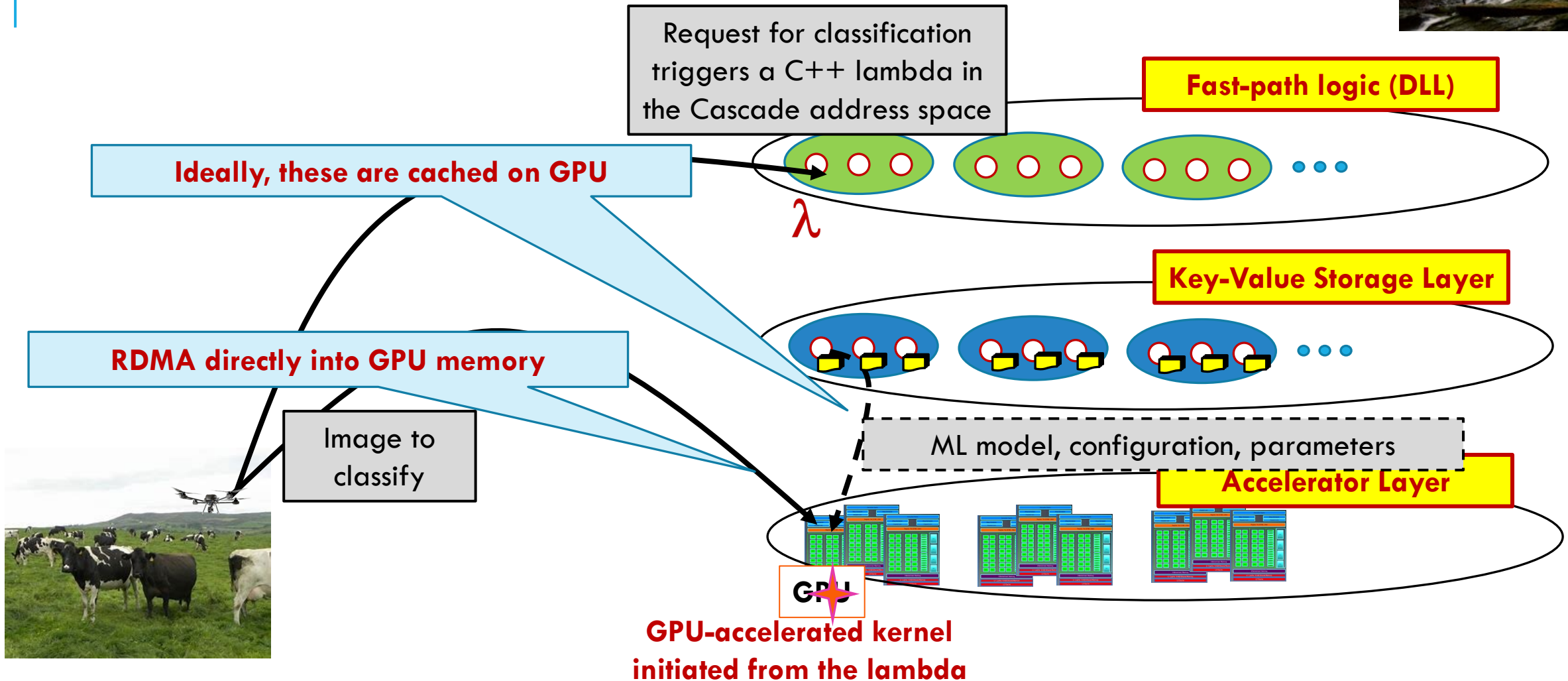


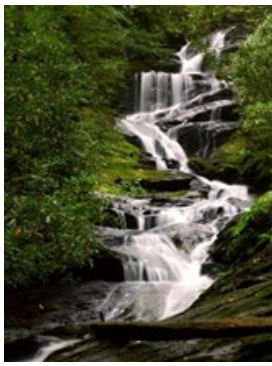
... BUT CAN ALSO HOST USER CODE





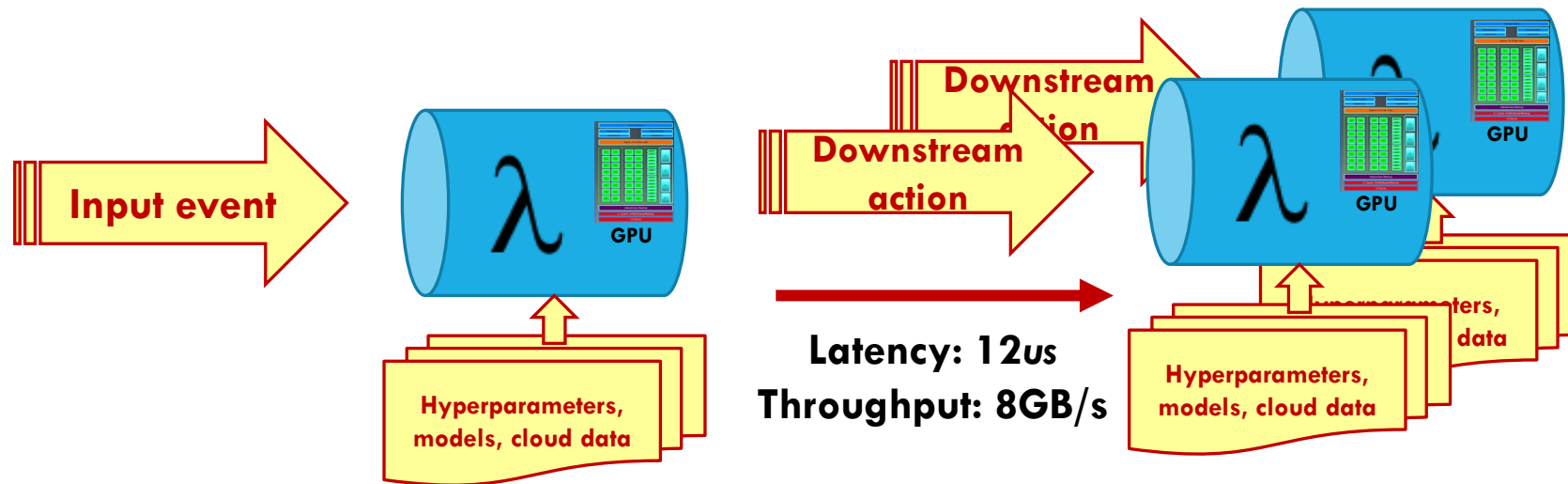
CASCADE: CUSTOMIZED SMART SERVICES



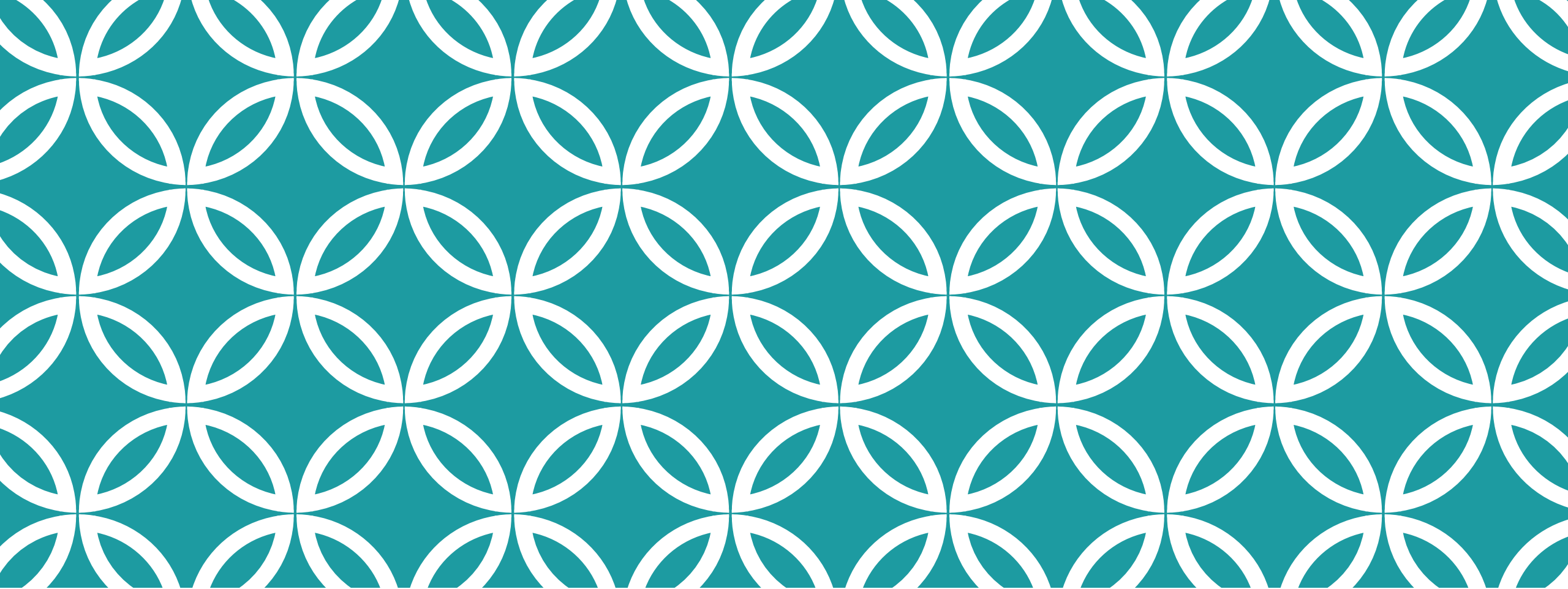


FAST-PATH PERFORMANCE

A simple federated ML pipeline



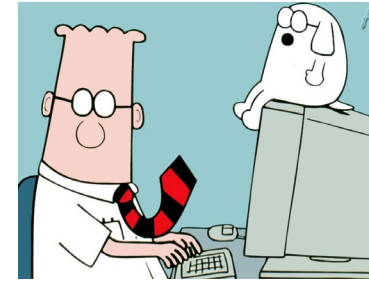
Cascade is close to ideal efficiency on our hardware and 100 to 10,000x faster than common options like Apache Flink



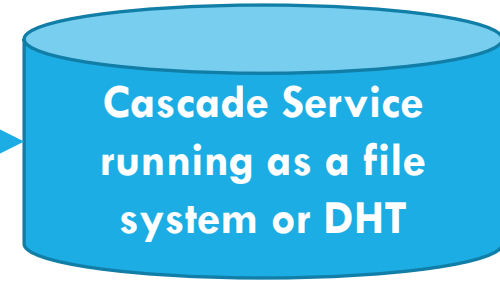
THE CASCADE MODEL

**Is it a service? A
library?**

CASCADE IS A SERVICE



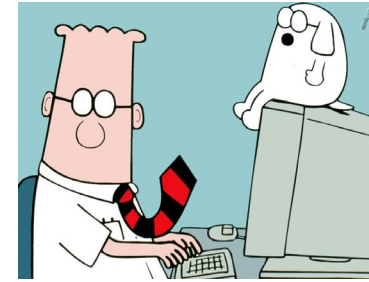
External Client



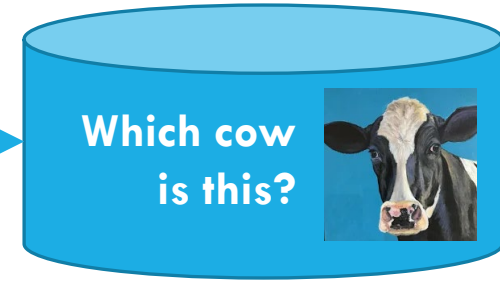
Cascade runs on a set of nodes (machines or VMs) where it controls some resources (cores, RDMA interfaces, GPUs/FPGAs, memory).

Users can build applications that access Cascade from “outside”. We call those “external clients”. They would think of Cascade purely as a key-value storage system, accelerated by RDMA.

EXTENSIBLE CASCADE IS AN SERVICE



External Client



But the highest speed is achieved by extending the Cascade service by adding logic (λ s) that will run inside Cascade.

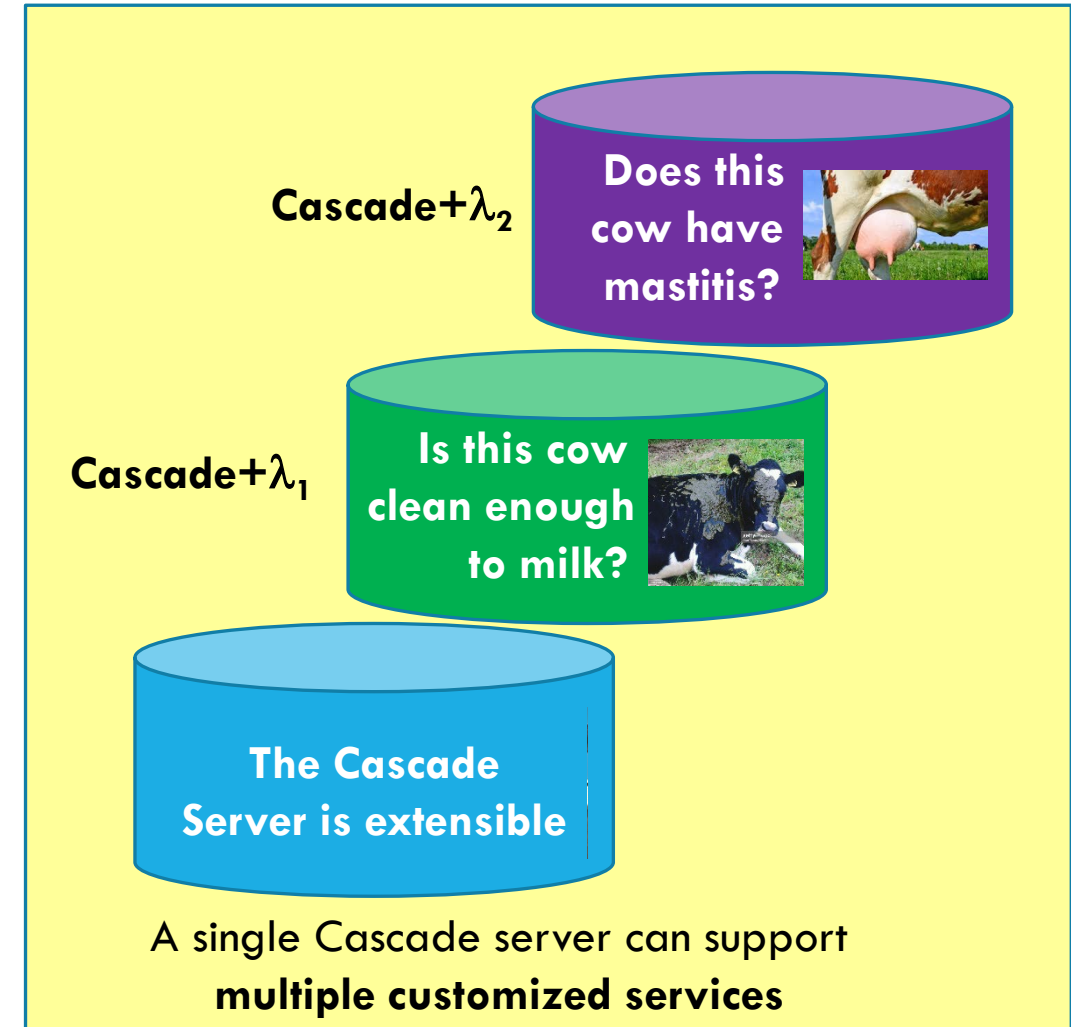
Used this way, Cascade enables creation of a customized service: a “smart” service ideal for your purposes.

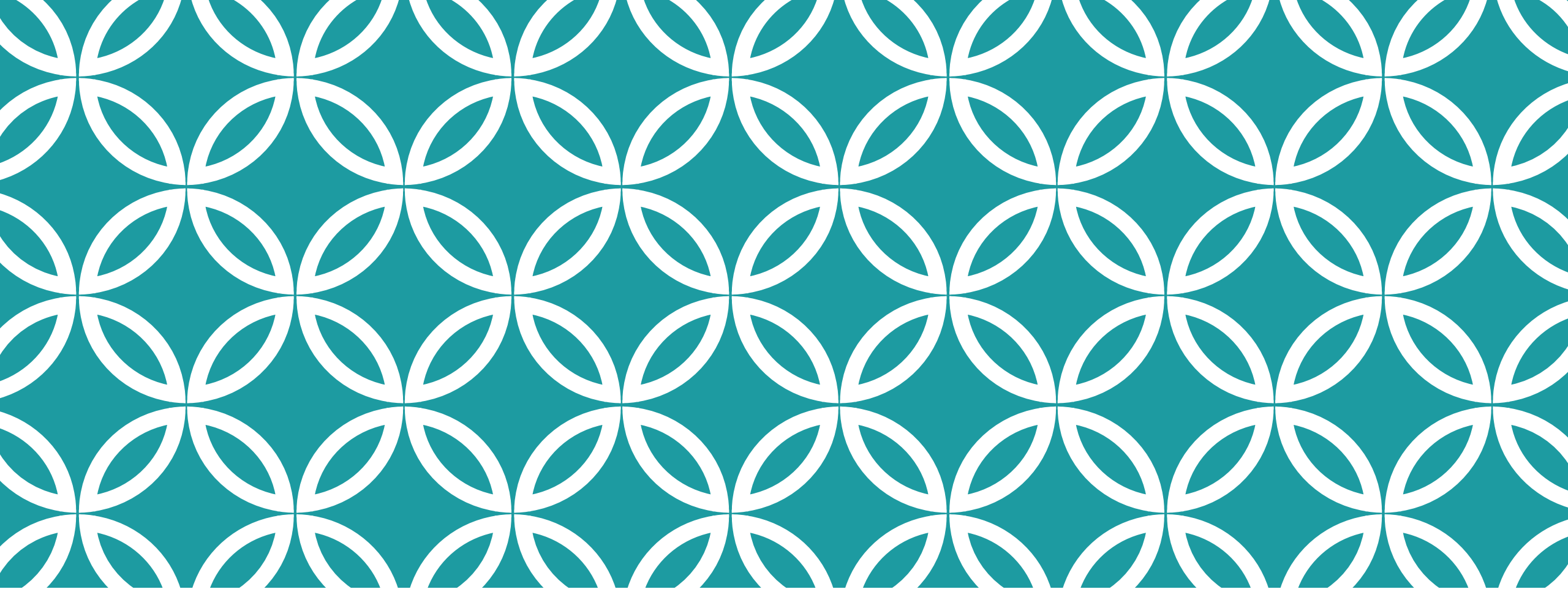
RECAP: EXTENSIBILITY (PAAS MODEL)

Cascade is one service

But when you supply customization
it acts like many specialized
services, one per application

So it becomes a platform for new
microservices, like these!





THE CASCADE STORAGE MODEL

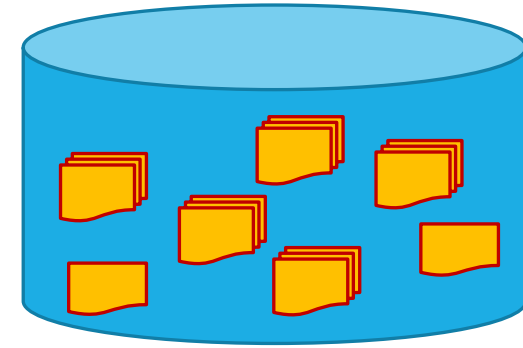
How should Cascade data be managed?

CASCADE IS A KEY-VALUE STORE

... but it also supports a file system API

... and the key-value API itself has some fancy options, beyond the basic **put/get/watch** we learned about in lectures 1-4

IN FACT CASCADE'S FILE SYSTEM IS AN EXAMPLE OF EXTENSIBILITY



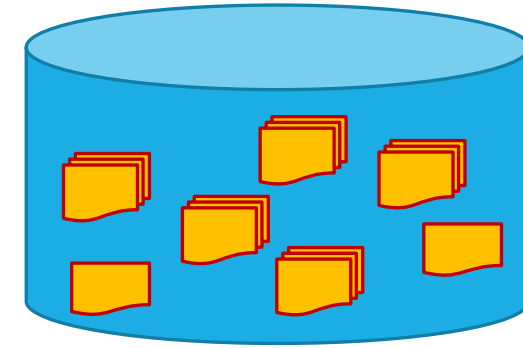
User sees scalable,
“private” object pools

Cascade's file system is implemented as a λ :

- Every object has a pathname.
- The file system extension supports normal file operations.
- You can access it just like any file system.

Similar to Ceph file system, yet just a few lines of code to translate user requests (via the FUSE library) into Cascade put/get/watch

UPDATES: CREATE A NEW FILE OR APPENDING TO AN EXISTING FILE



User sees scalable,
“private” object pools

Same model used in cloud file systems like HDFS and Ceph.

- Updates are log appends using Paxos. Each object has a log of versions that evolved over time.
- Reads run on the *stable prefix* of the log.

VERSIONED UPDATES

Each time you write to an object, Cascade creates a new version.



If you read an object, modify it, and then write it back, you can tell Cascade which version you modified. The **put** will double check to be sure that this is still the current version before replacing it, and otherwise returns an error (then you can loop)

VERSIONED/TEMPORAL QUERIES

Accessed via **get**, but specify the version # or time desired

In the volatile case, Cascade only keeps the most recent version. With persistent objects, Cascade keeps a log of past versions.

- By default, applications see the most current version
- Indexed access allows the application to query any version (by version number or time), or fetch any data range.

VERSIONED OBJECTS

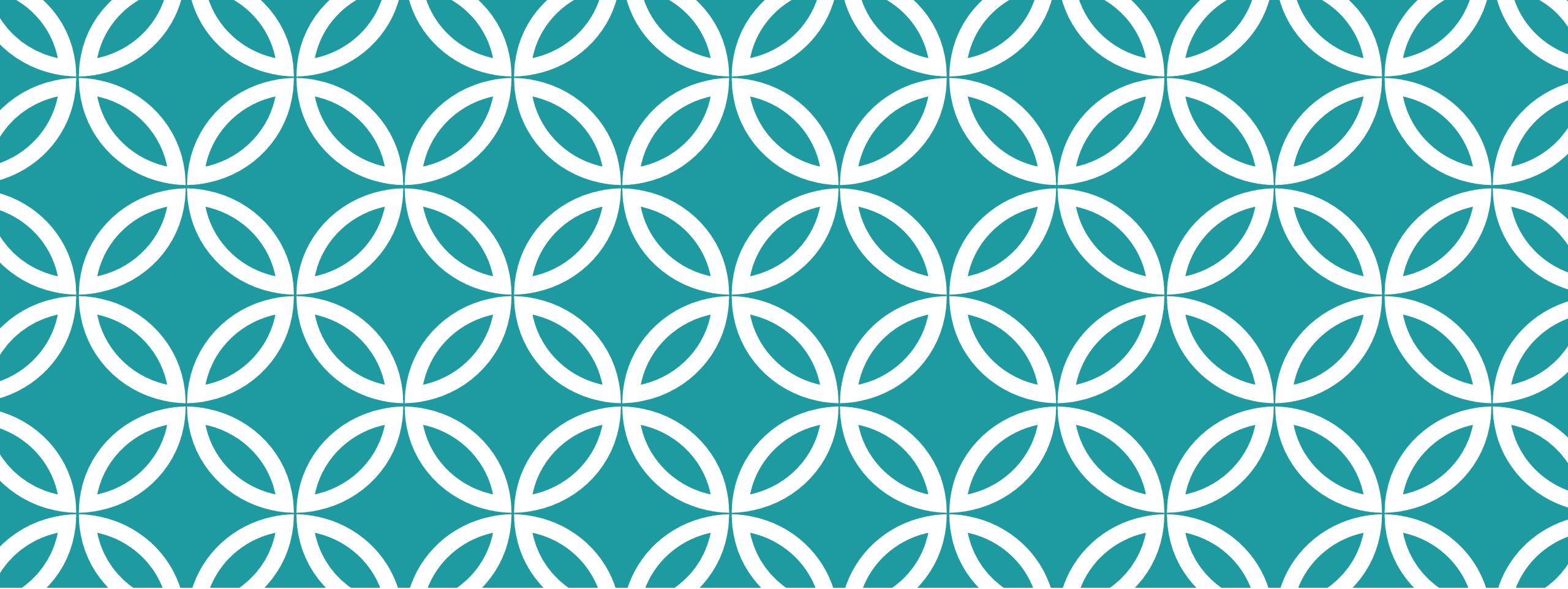


Versioning enables intelligence over a sequence of events

- *key*: The object store *always* tracks information on a per-object basis
- *version-number*: Just an integer
- *time*: If the object itself lacks a timestamp, we just use “platform” time.

Now **get** can lookup most current version, or a specific one, even by time.

The object store is optimized to leverage non-volatile memory hardware.



THE CASCADE COMPUTE MODEL

**Virtual synchrony,
Atomic multicast, Paxos,
consistent cuts...**

THE CENTRAL PUZZLE

The very fastest data paths require compilation, ideally in languages like C++.

But we want Cascade to run as a *service*, so it would often already be running when a new user comes along and wishes to create and launch some completely new service.

How can we extend a running system? Actually... it isn't so hard

FIRST QUESTION: WHAT'S IN A λ ?



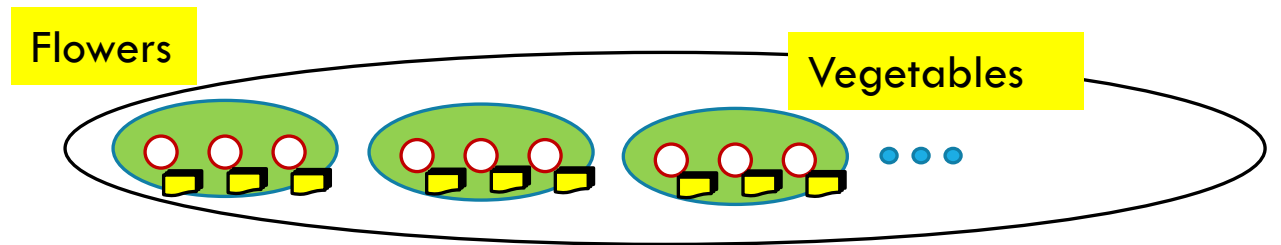
We support many languages. Native APIs are Python with various packages (including LINQ) and C++ with LINQ.

Code is concise – LINQ pioneered a style that mixes “kernel” invocations with embedded SQL. Maps cleanly to GPU, FPGAs.

Cascade manages GPUs and can cache data in GPU memory.

HOW CAN A KEY-VALUE STORE “BE” A CLASSIFIER SERVICE OR AN ANALYTIC SERVICE?

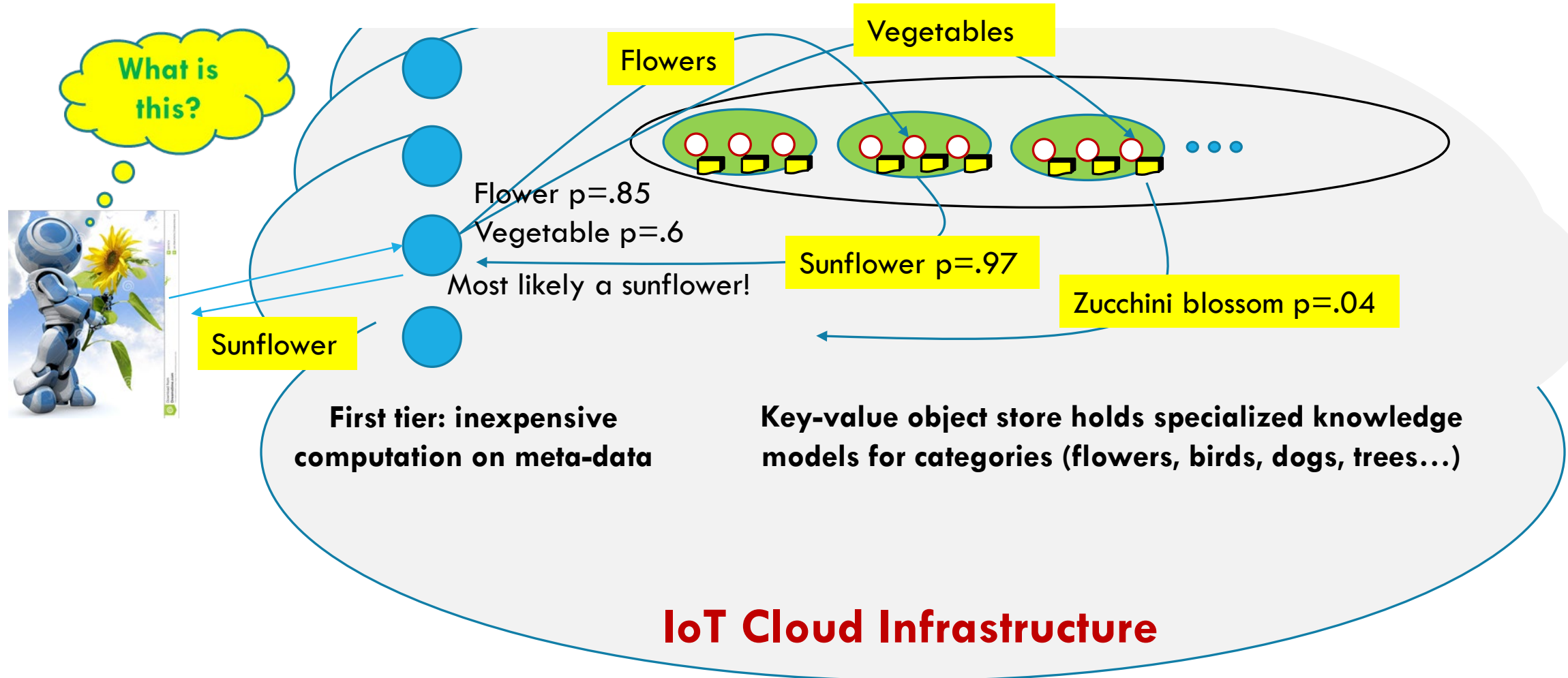
We run Cascade on a set of nodes. Here we see nine nodes in three shards.



A shard identically replicates (key,value) tuples, using Paxos.

Here, an object with the key “Flowers” was stored in shard 0. “Vegetables” ended up in shard 2.

... ENABLING THIS KIND OF SOLUTION



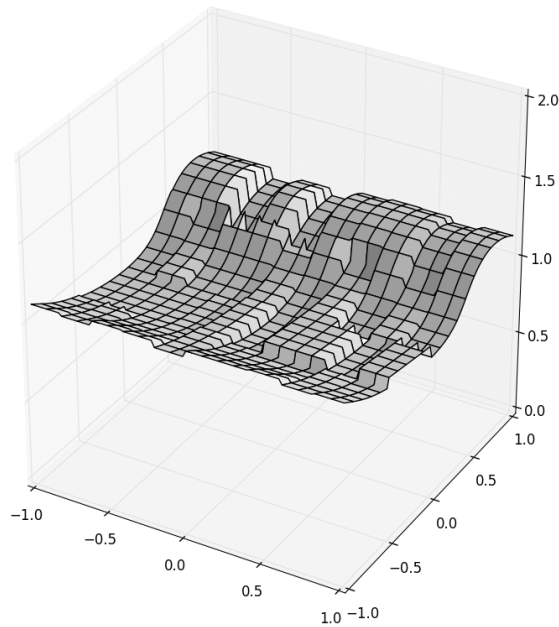
CASCADE QUERIES SEE CONSISTENT CUTS

Suppose a cooperative AI is triggered by event ε at time τ . We run all the lambdas triggered by ε along a consistent cut “optimally close” to time τ (and selected deterministically).

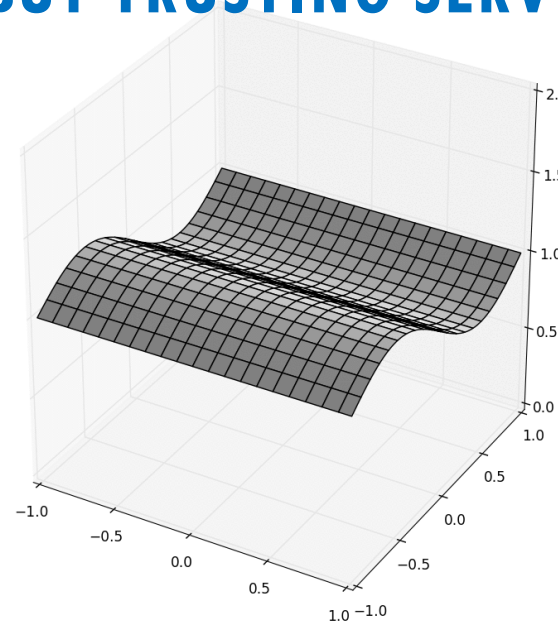
Effect: The lambda won't see platform-induced inconsistencies.

VISUALIZATION OF CASCADE CONSISTENCY

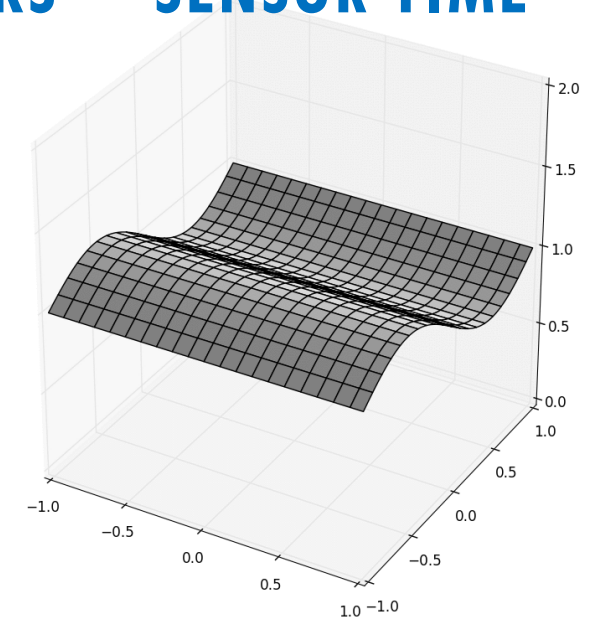
HDFS



CASCADE USING CONSISTENT CUTS BUT TRUSTING SERVER CLOCKS



CASCADE WITH SENSOR TIME



Cascade consistent cuts + GPS-timestamped sensor data result in clean input to the D-AI algorithm (in this case, a simple visualization)

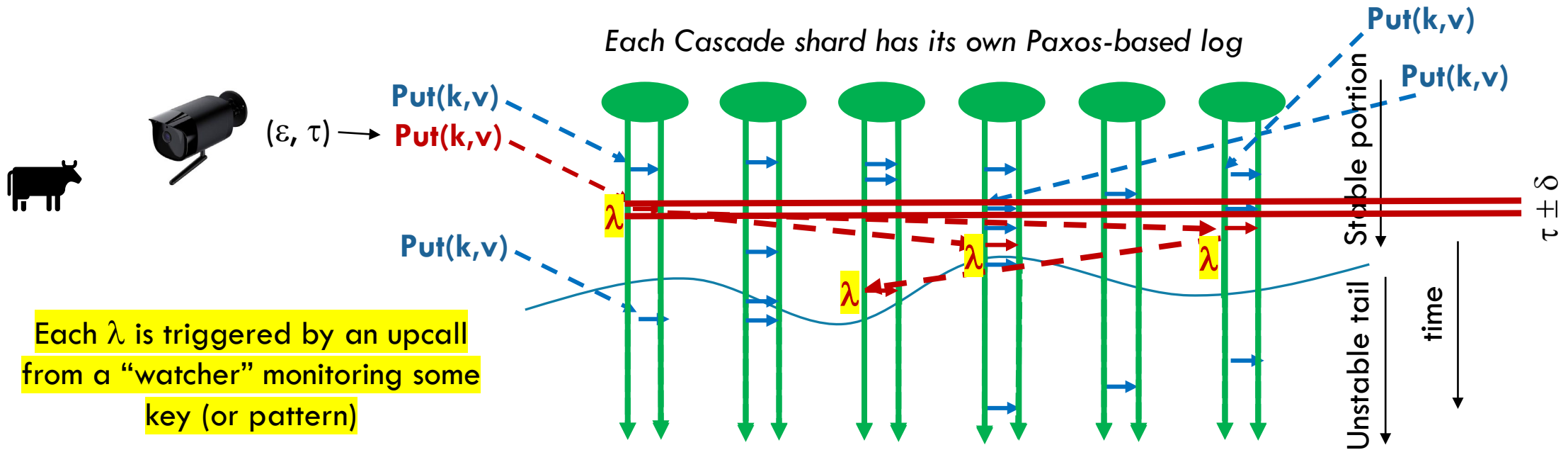
CENTRAL CONCEPTUAL INSIGHT

One event may trigger many lambdas.

These lambdas may need to run on multiple nodes... yet will share the same temporal index (τ from the trigger event ε).

A Cascade query always sees a “consistent state snapshot.”

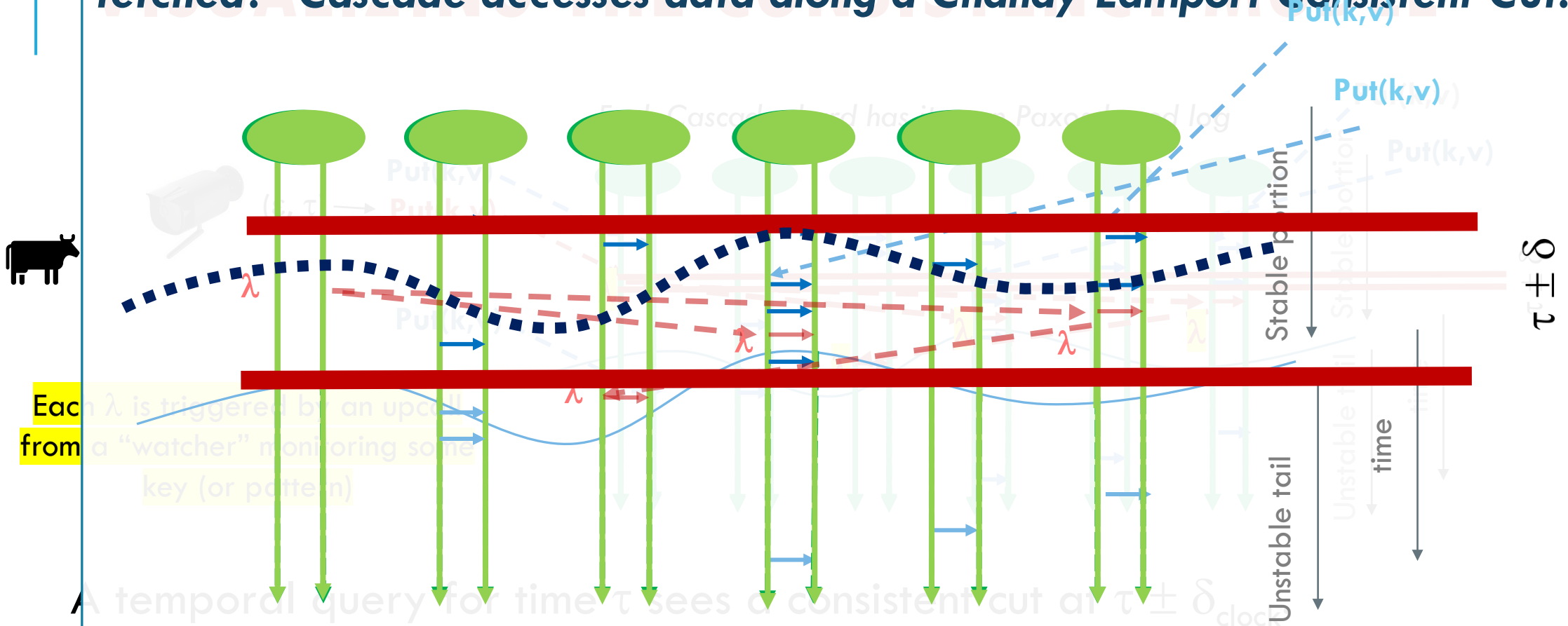
VISUALIZING THE CONSISTENCY MODEL



A temporal query for time τ sees a consistent cut at $\tau \pm \delta_{\text{clock}}$.

Queries to unstable data must wait, but updates are stable within 50us.

Even if δ is small, there is a choice to make: Which versions will be fetched? Cascade accesses data along a Chandy-Lamport Consistent Cut.

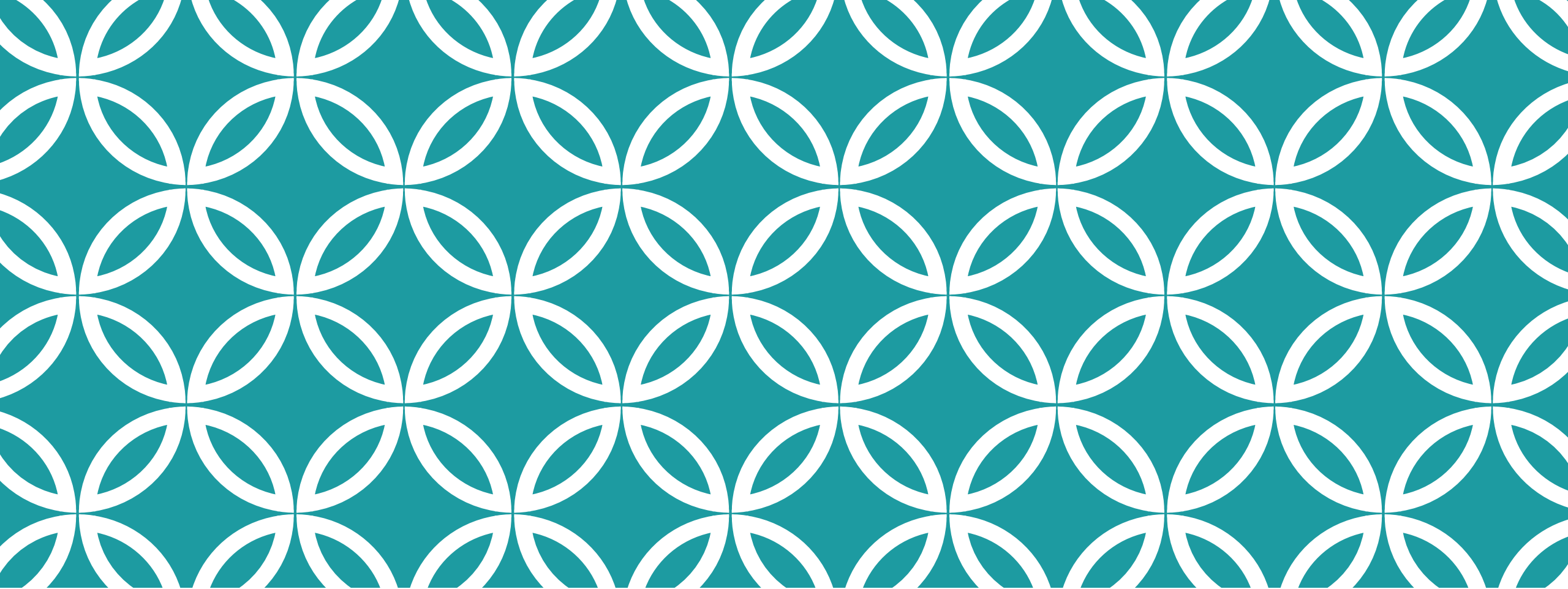


GROUPING OBJECTS

Often a lambda will need to access several objects that should ideally all have their own keys, yet you want them grouped on the same shard.

For this, Cascade supports “affinity grouping”. Each object has a second key, used for placement.

Even if A and B have different keys — “names” — they will be stored on the same shard if you assign them the same affinity key.



WHAT DOES A REAL APPLICATION LOOK LIKE TODAY?

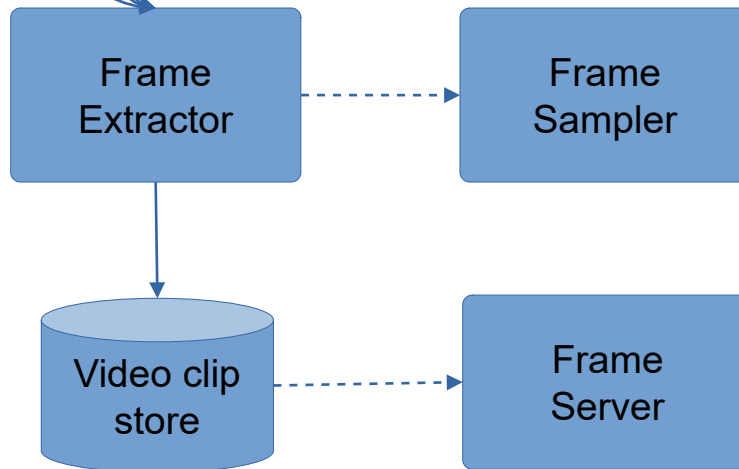
**Example, courtesy of Weijia,
Alicia and Thompson**

DAIRY IMAGE PIPELINE: FRONT END

Dairy Farm

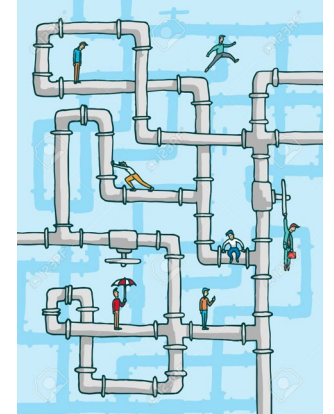


The Farm Server (IoT Edge)

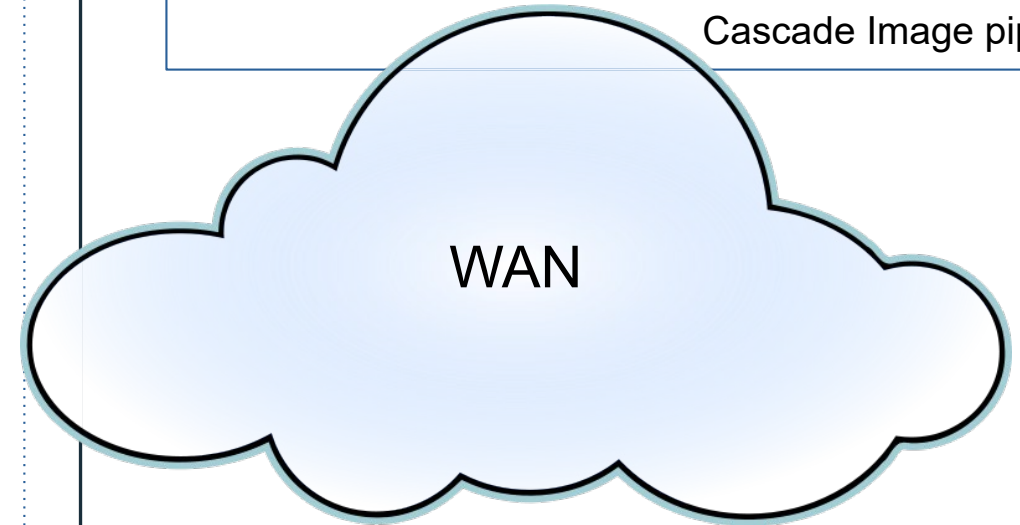


Data Center

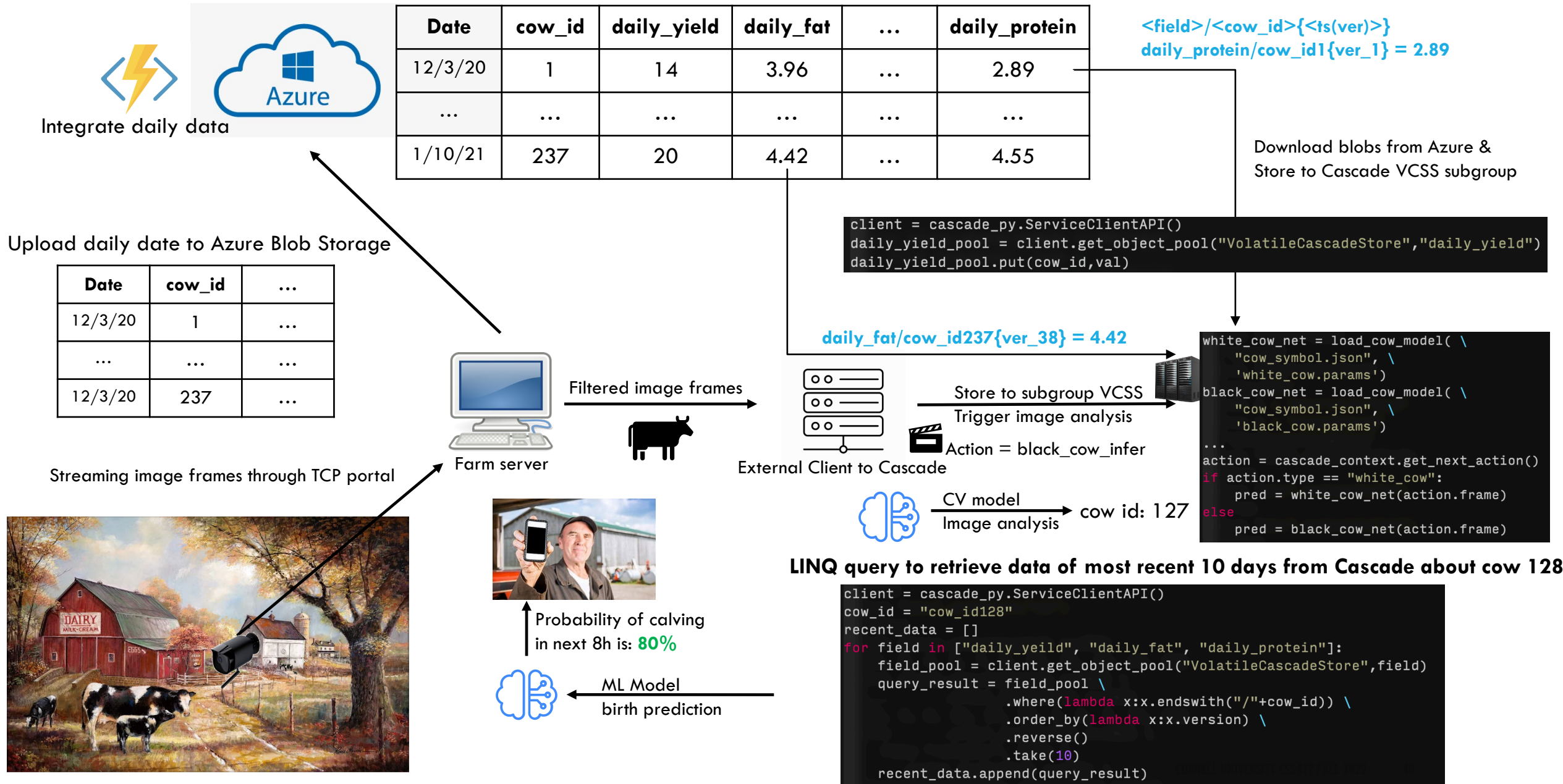
Image Pipeline
Front End
(As an external client)



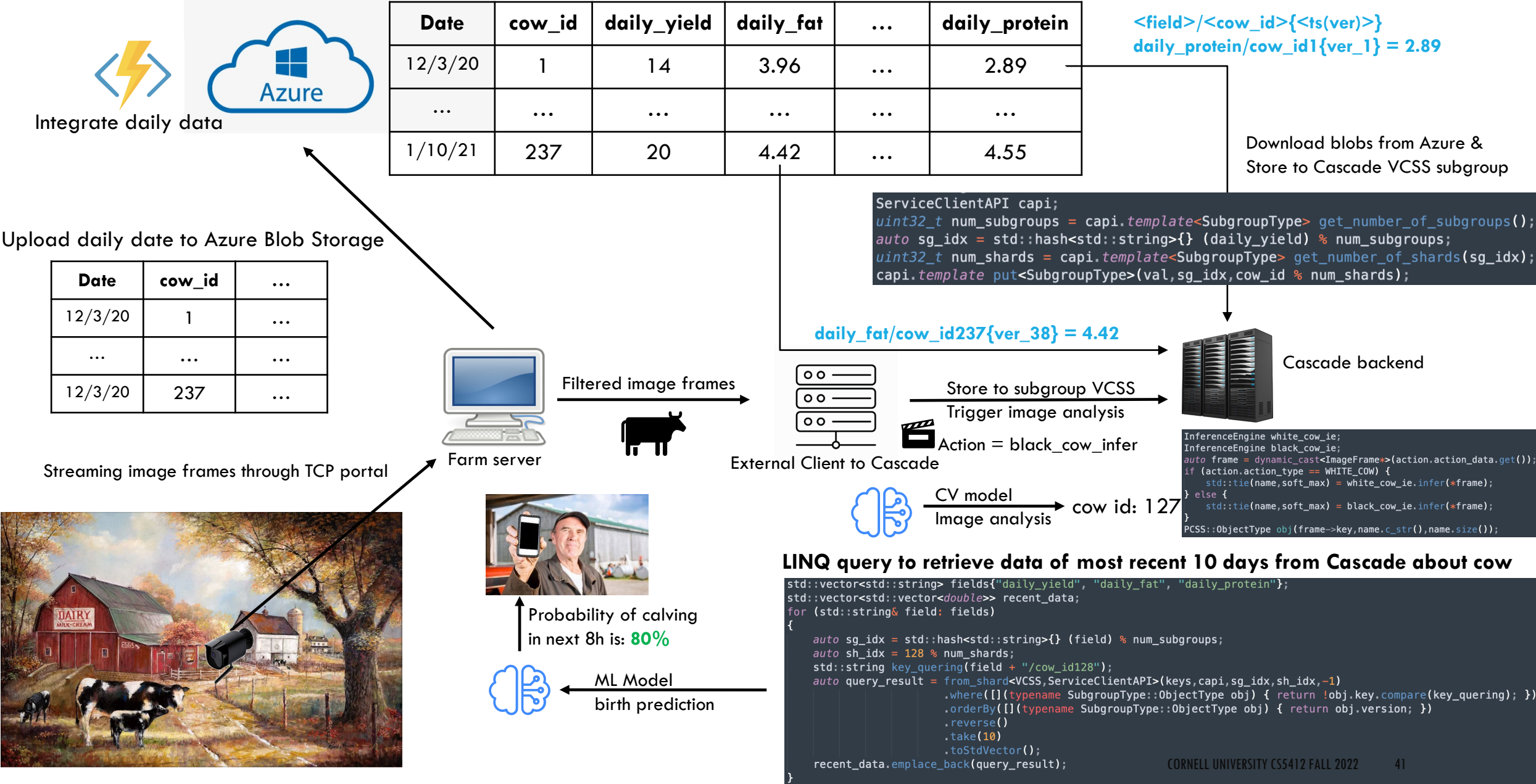
Cascade Image pipeline



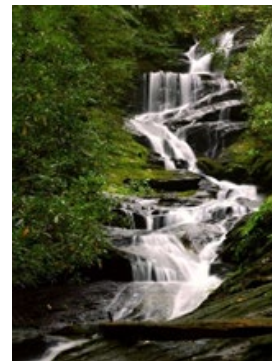
... DETAILED VERSION (PyLINQ ON MSFT AZURE)



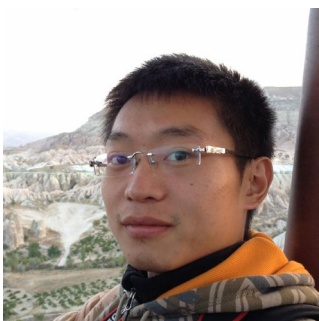
C++ IS SIMILAR (BUT MORE EFFICIENT)



THE CASCADE AND DERECHO TEAM



Weijia Song



Alicia Yang

Ken Birman



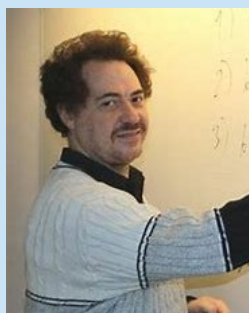
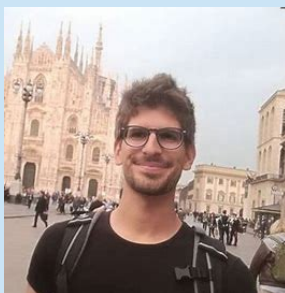
Sagar Jha

Lorenzo Rosa



Mae Milano
(post-doc at Berkeley)

Andrea Merlina



Roman Vitenberg

Cornell undergrads:

Aahil Awatramani, Ben Posnick, Max
Charlamb, Archishman Sravankumar, Aaron
Weiss, Peter Zheng



Thompson Liu



Edward Tremel
(faculty at Augusta Univ.)