# *EdgStr*:
# Automating **Client-Cloud** to **Client-Edge-Cloud** Transformation



Principle Engineer
Software Engineering Team
Samsung Research

Kijin An



Professor
CS Department
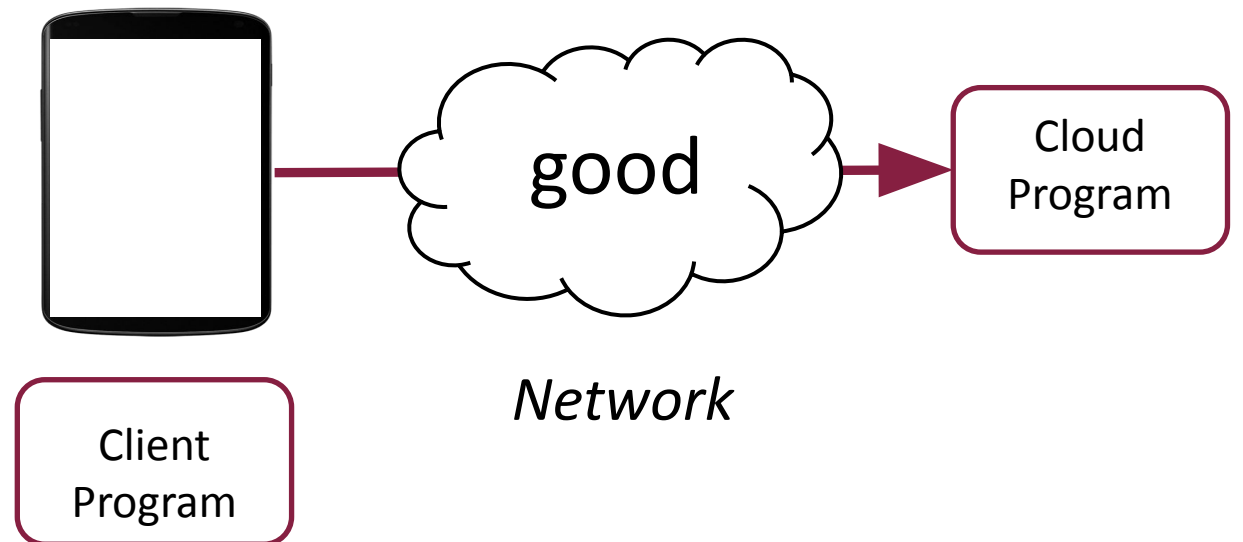Virginia Tech

Eli Tilevich

VIRGINIA TECH.

# Presentation Outline

- Motivation

    ○ State of distributed software and vision

    ○ Example application

- Approach for automating transformation to Client-'Edge'-Cloud

- Reference Implementation: EdgStr
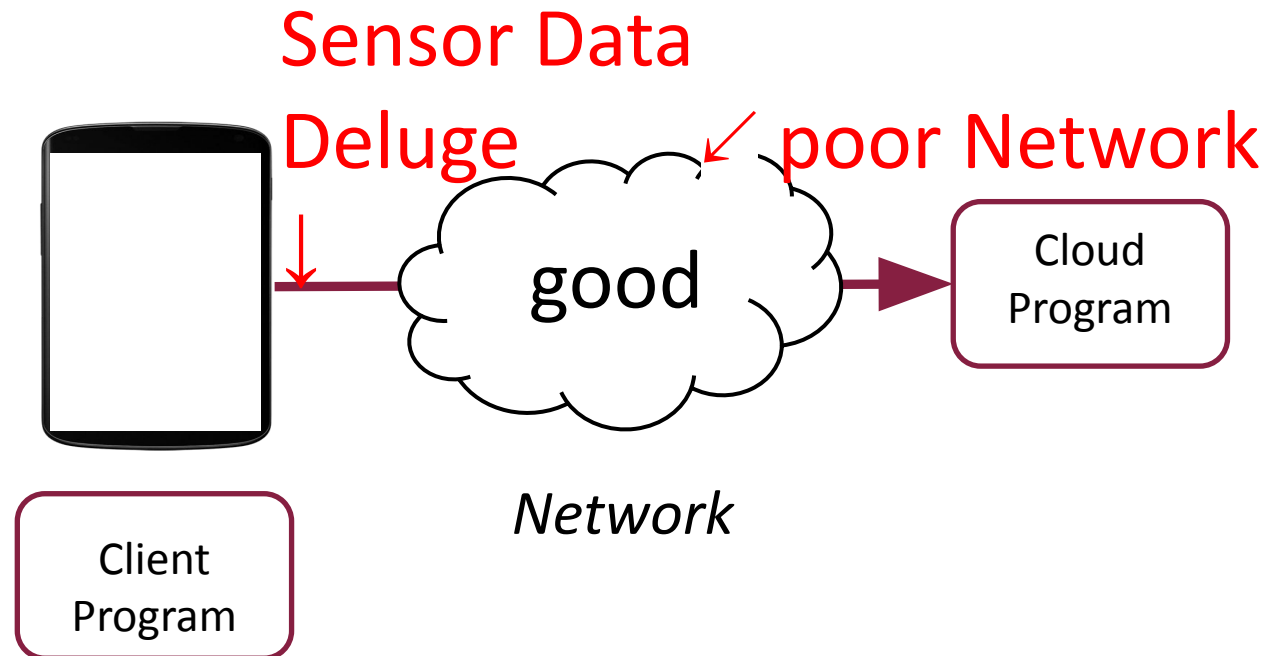
- Evaluation

- Conclusions

# Client-Cloud Architecture (2-tier)

- Cloud-Client predominant
  - Cloud Infrastructure: Powerful
  - Network: Fast

- Conventional 2-tier no longer meets performance and resource utilization requirements of modern apps

good

Cloud Program

Network

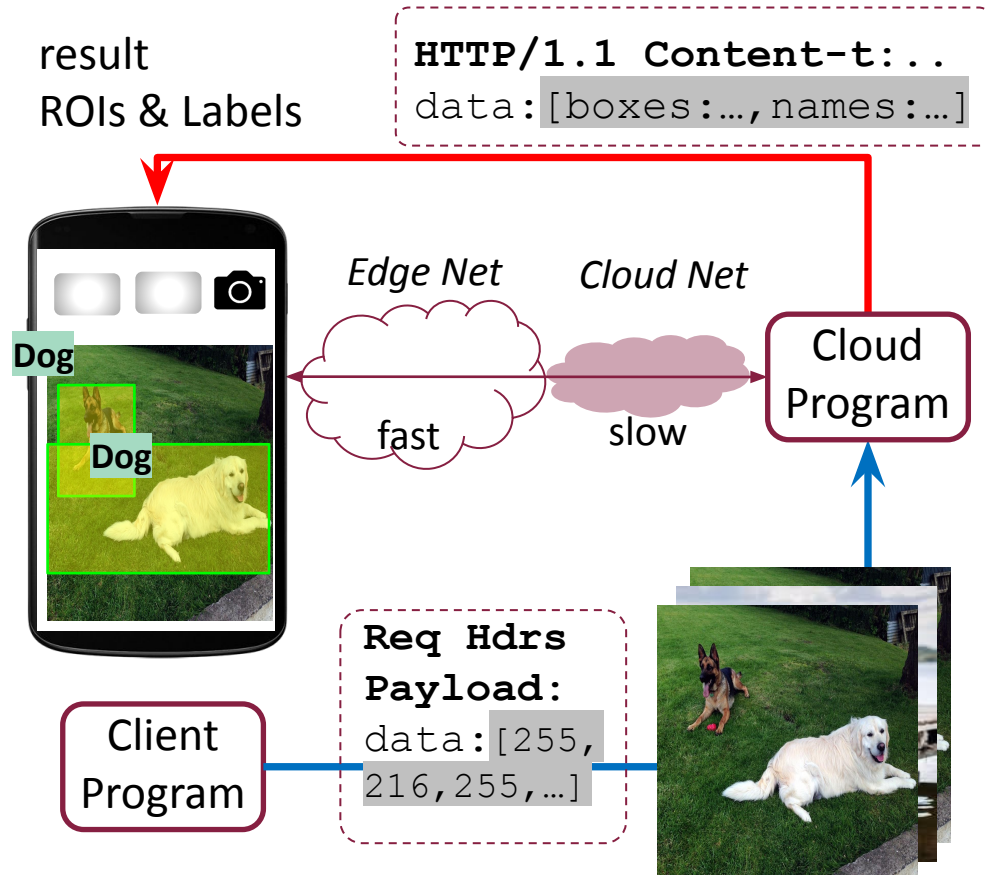Client Program

# Client-Cloud Architecture (2-tier)

- Cloud-Client
  - Cloud Infrastructure: Powerful
  - Network:Fast

- What if:
  - Network: Slow & unreliable
  - Sensor Data: Massive ("sensor deluge")
  - ☐ Increased Latency

**Sensor Data Deluge**

**poor Network**

good

Cloud Program

Client Program

*Network*

4

# Motivating Example
(firebase-objdet)

- Client-Cloud program (/predict, detect objects in the cloud)

result
ROIs & Labels

**HTTP/1.1 Content-t:..**
data:[boxes:…,names:…]

*Edge Net*    *Cloud Net*

Dog

Dog

Cloud
Program

fast    slow

# Sensor Data Deluge!

**Req Hdrs**
**Payload:**
data:[255,
216,255,…]

Client
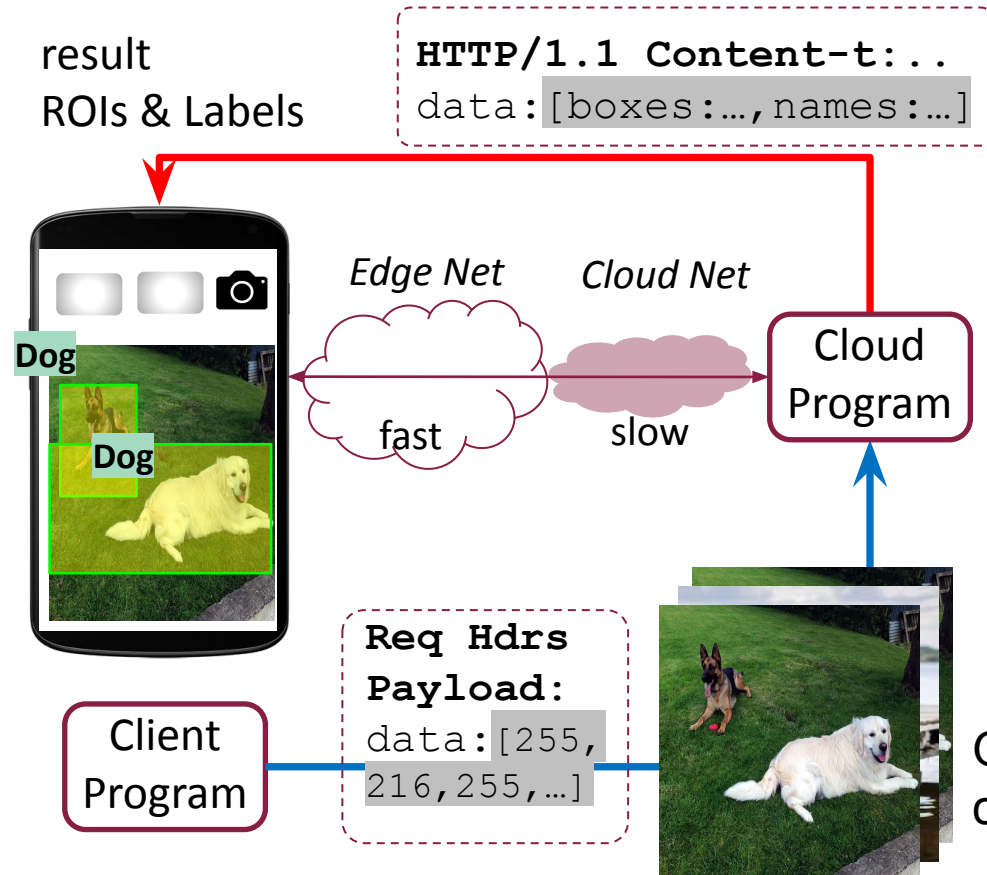Program

Galaxy S24 Ultra or iPhone 15 Pro
can capture a photo of **12MBytes**

5

# Motivating Example
(firebase-objdet)

- Client-Cloud program



result
ROIs & Labels

```
HTTP/1.1 Content-t:..
data:[boxes:…,names:…]
```

Dog

Dog

*Edge Net*    *Cloud Net*

Cloud Program

fast    slow

Client Program

```
Req Hdrs
Payload:
data:[255,
216,255,…]
```

**RTT across different/same continents are different from**
:An Order of magnitude between them!

Cloud Program
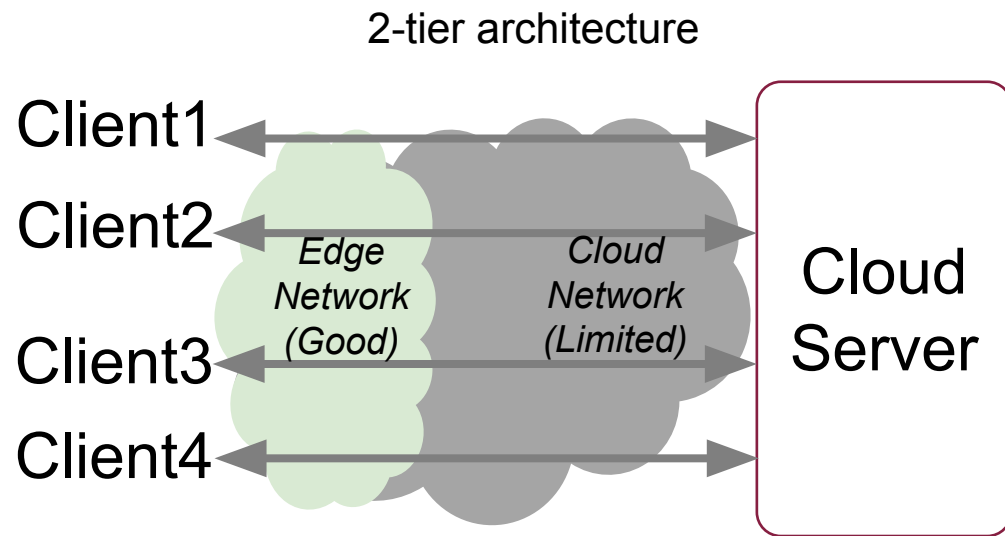
Cloud Program

Dog
Dog

[Installed Cloud Programs differently on *Heroku* platform]

Galaxy S24 Ultra or iPhone 15 Pro can capture a photo of **12MBytes**

6

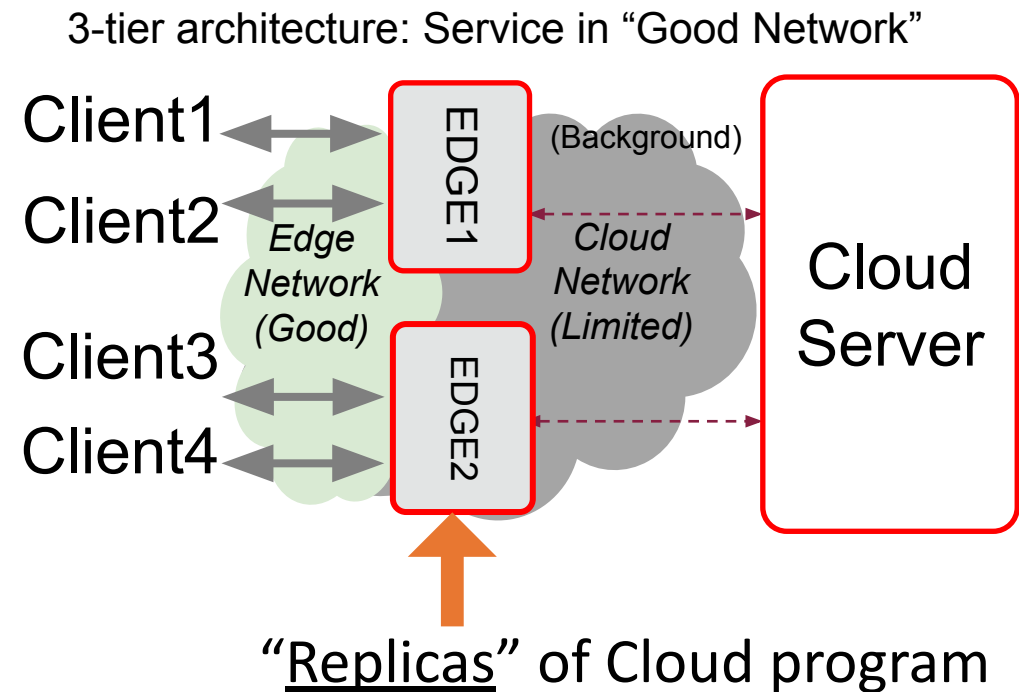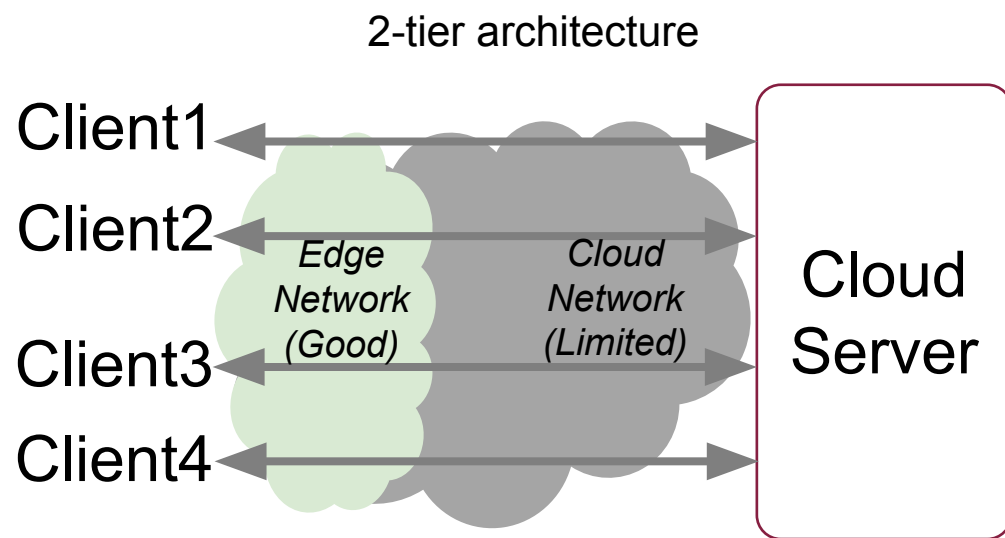# Transforming 2-tier into 3-tier architecture

- Edge-based processing benefits

2-tier architecture

Client1

Client2

Client3

Client4

Edge Network (Good)

Cloud Network (Limited)

Cloud Server

# Transforming 2-tier into 3-tier architecture

- Benefit from edge-based processing

2-tier architecture

Client1

Client2

Client3

Client4

*Edge Network (Good)*

*Cloud Network (Limited)*

Cloud Server

3-tier architecture: Service in "Good Network"

Client1

Client2

EDGE1

(Background)

*Edge Network (Good)*

*Cloud Network (Limited)*

Cloud Server

Client3

Client4

EDGE2

"Replicas" of Cloud program

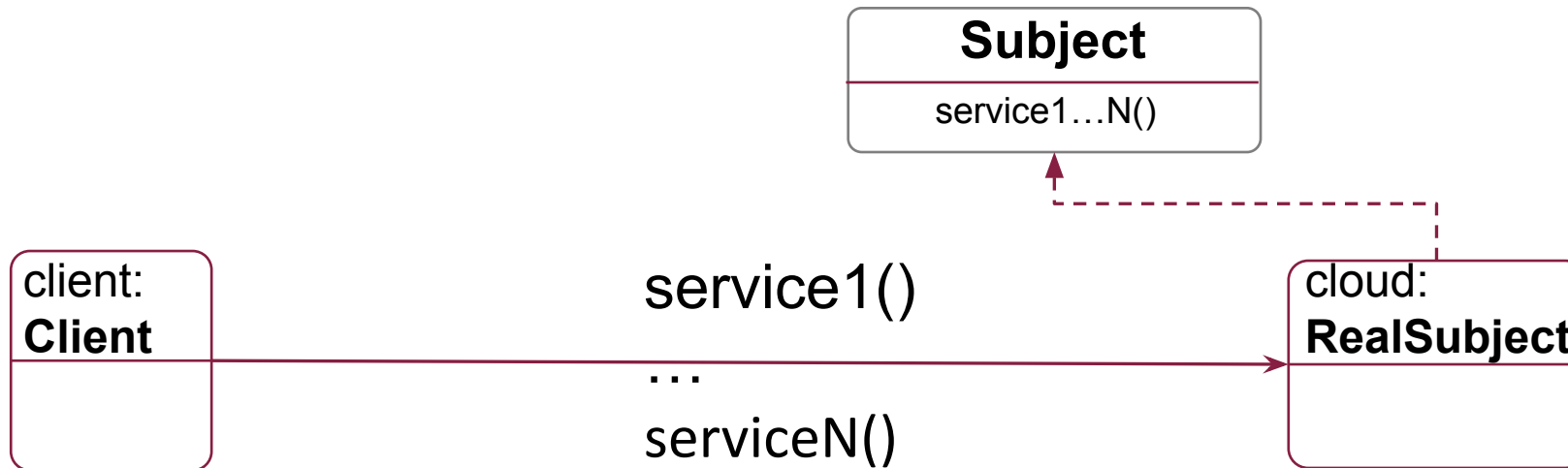- We understand the benefits, but how to automate the transformation?

# Approach Overview

- Replicate a cloud-based service on edge devices
- Select the portion functionality to replicate that improves performance
- Provide eventual consistency by relying on CRDT
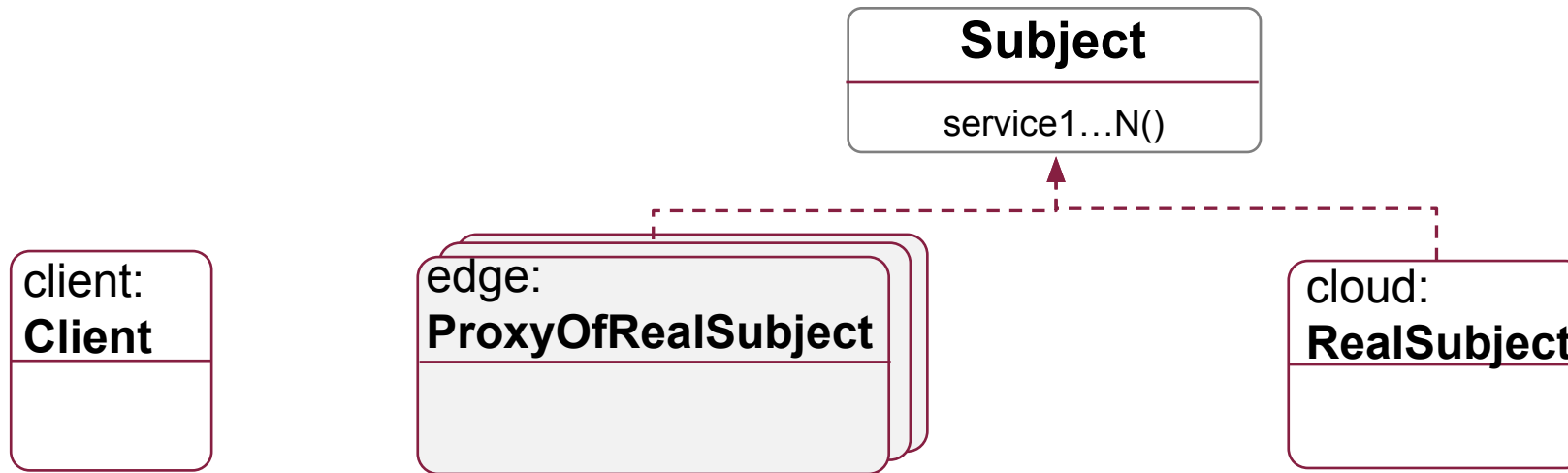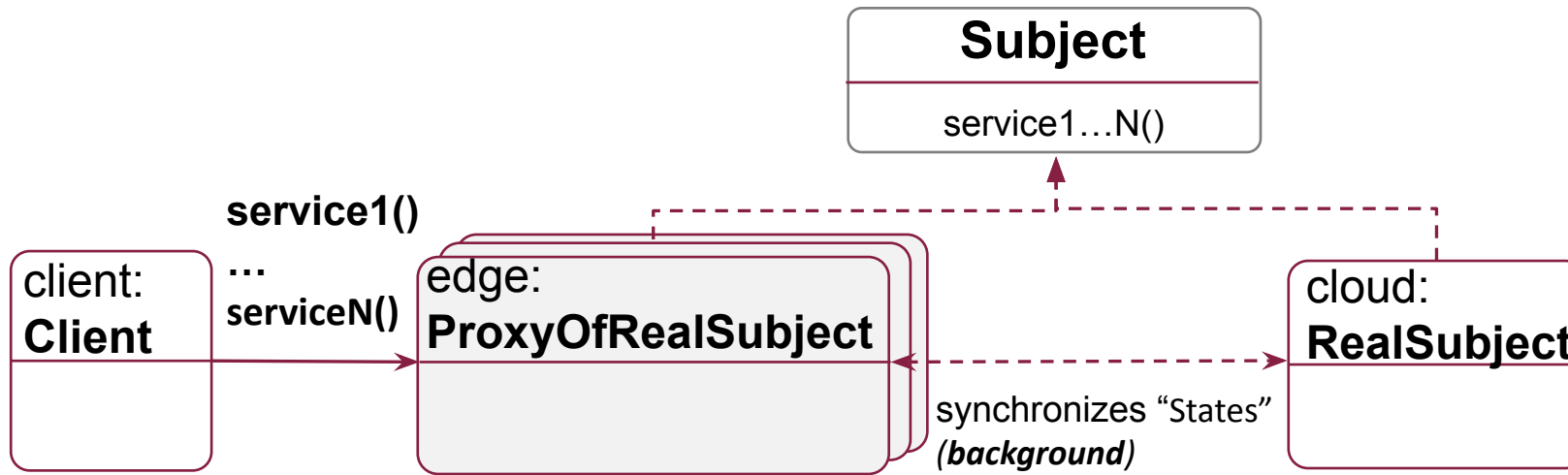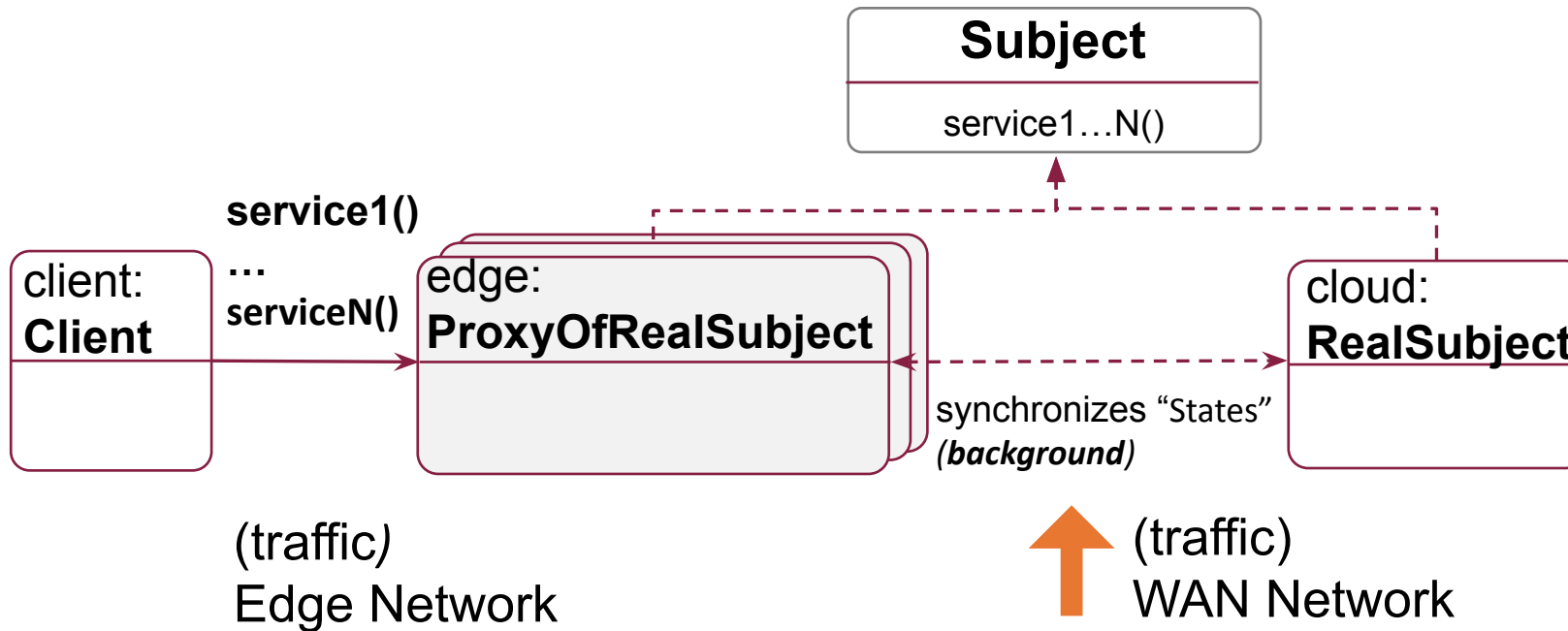- Load-balance to a cluster of edge devices for scalability and throughput

# Edge Processing via Tailed Proxy Pattern

- Proxy Pattern: Client makes request to Proxy (edge replicas)

# Edge Processing via Tailed Proxy Pattern

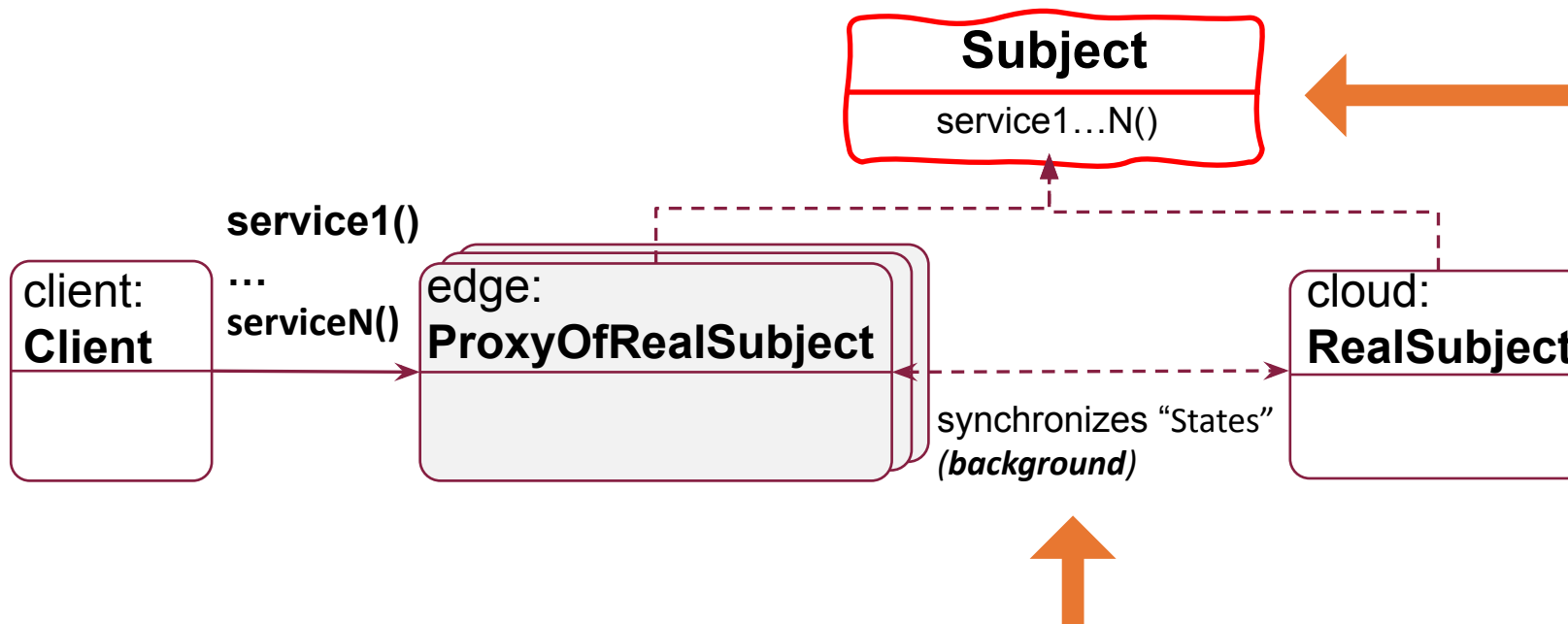- Proxy Pattern: Client makes request to Proxy (edge replicas)

# Edge Processing via Tailed Proxy Pattern

- Proxy Pattern: Client makes request to Proxy (edge replicas)



**service1()**
...
**serviceN()**

| **Subject** |
| --- |
| service1…N() |

client:
**Client**

edge:
**ProxyOfRealSubject**

cloud:
**RealSubject**

synchronizes "States"
*(background)*

# Edge Processing via Tailed Proxy Pattern

- Proxy Pattern: Client makes request to Proxy (edge replicas)



**Subject**

service1…N()

**service1()**
…
**serviceN()**

client:
**Client**

edge:
**ProxyOfRealSubject**

cloud:
**RealSubject**

synchronizes "States"
*(background)*

(traffic*)*
Edge Network

(traffic)
WAN Network

relaxed consistency:
synchronized in a background process without
interfering with main functionality

# Edge Processing via Tailed Proxy Pattern

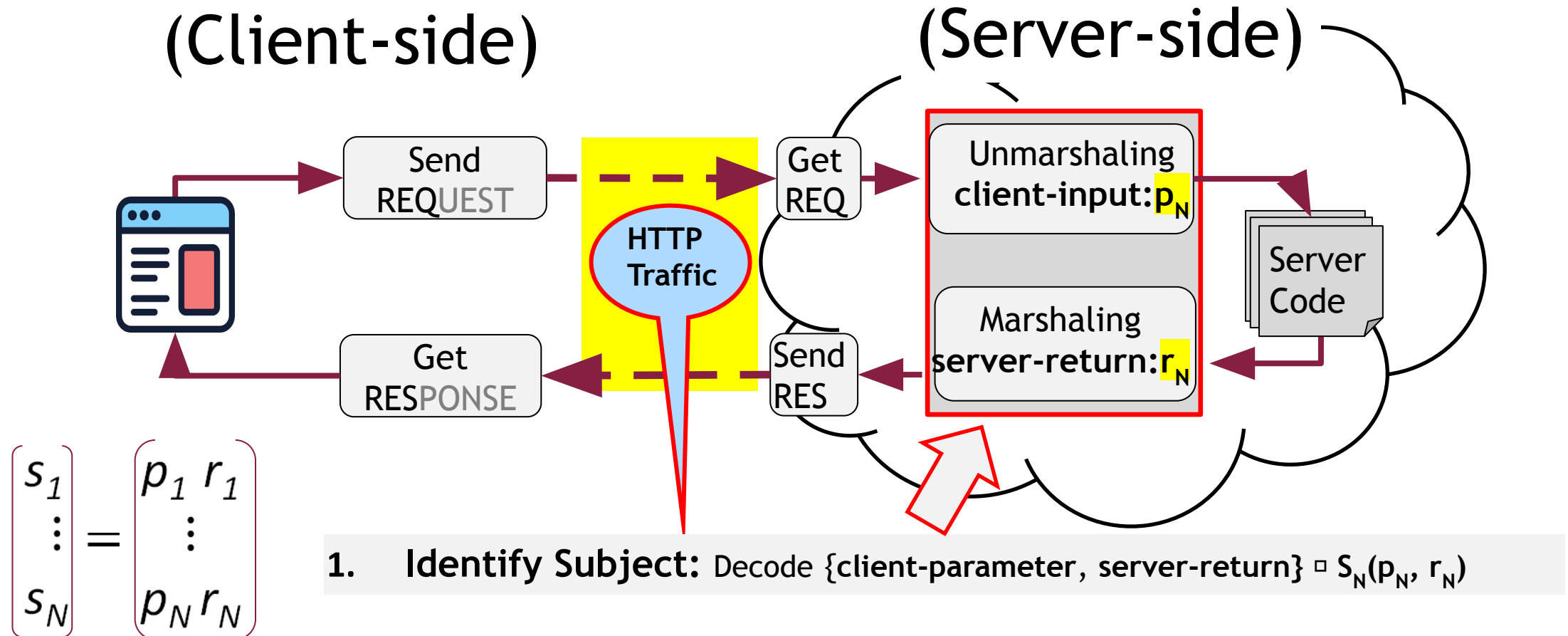- Proxy Pattern: Client makes request to Proxy (edge replicas)

**Subject**

service1…N()

**Program Analysis & Transformation:**
How to identify and extract required subject functionalities in Cloud program?
: carefully choosing to benefit from edge based processing!

**service1()**
...
**serviceN()**

client:
**Client**

edge:
**ProxyOfRealSubject**

cloud:
**RealSubject**

synchronizes "States"
*(background)*

**relaxed consistency:**
synchronized in a background process without interfering with main processing

14

# EdgStr: Automated Transformation

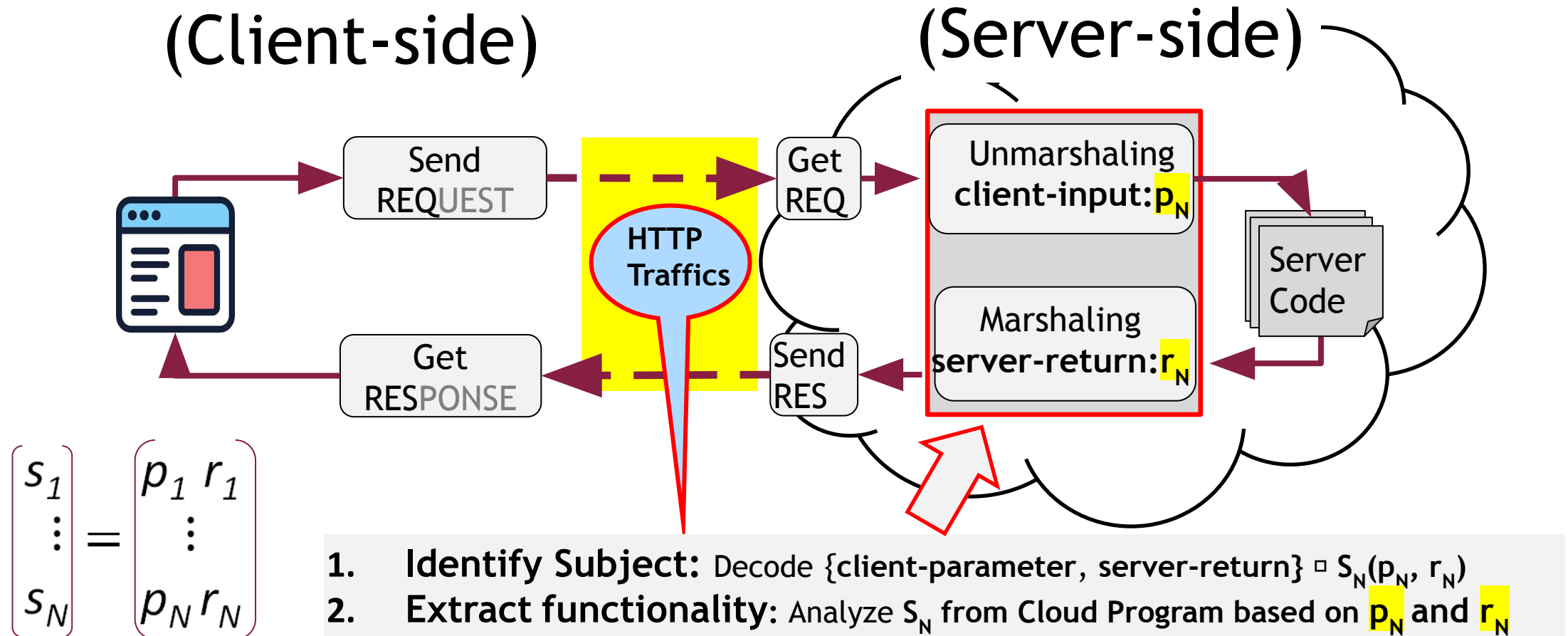Identify and extract required subject functionalities in Cloud program

- Identify Subject $s_1$, $s_2$,...$s_N$ from by capturing HTTP traffic

(Client-side)                    (Server-side)



$$\begin{bmatrix} s_1 \\ \vdots \\ s_N \end{bmatrix} = \begin{bmatrix} p_1\ r_1 \\ \vdots \\ p_N\ r_N \end{bmatrix}$$

1. **Identify Subject:** Decode {client-parameter, server-return} ⊡ $S_N(p_N, r_N)$

# EdgStr: Automated Transformation

Identify and extract required subject functionalities in Cloud program

- **Extracting** "functionality" from Cloud program:



(Client-side)     (Server-side)

$$\begin{bmatrix} s_1 \\ \vdots \\ s_N \end{bmatrix} = \begin{bmatrix} p_1 \ r_1 \\ \vdots \\ p_N \ r_N \end{bmatrix}$$

1. **Identify Subject:** Decode {client-parameter, server-return} ▫ $S_N(p_N, r_N)$
2. **Extract functionality:** Analyze $S_N$ from Cloud Program based on $p_N$ and $r_N$

# EdgStr: Overall Process

- Dynamic and Static Analysis for Cloud Program

# EdgStr: Overall Process

$$\textbf{ExtractedStmts}\,(s_n,\,V^{u_{id}}_{unMar},\,V^{u_{id}}_{Mar}) \leftarrow \;;;\,V^{u_{id}}_{unMar}:rN,\,V^{u_{id}}_{unMar}:p_N$$

$$\left(\textbf{Stmt\_Dep}(s_n,\,s_1)\wedge\textbf{Marshal}(s_1,\,v_{Mar},\,V^{u_{id}}_{Mar})\right)$$

$$\wedge\left(\neg\textbf{Stmt\_Dep}(s_n,\,stmt_2)\wedge\textbf{UnMarshal}(s_1,\,v_{unMar},\,V^{u_{id}}_{unMar})\right)$$

- Dynamic and Static Analysis for Cloud Program

# EdgStr: Overall Process

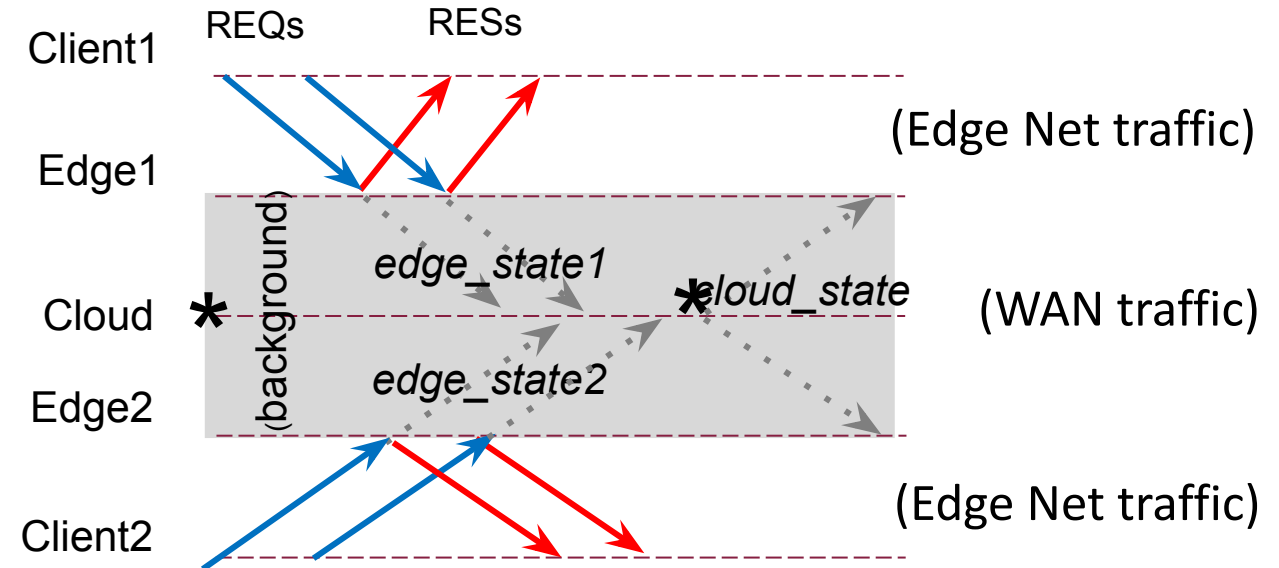- Dynamic and Static Analysis for Cloud Program

# States are synchronized: Between Cloud and Edge Replicas

- Eventually Consistency Sync. with CRDT for read or write operations across edge replicas and original cloud



Original
Client-Cloud app

Transformed
Cloud-Edge-Client

# Evaluation

- 7 open-source distributed apps (42 remote services)

- Edge Node Setup: RPI-3s and RPI-4s

- **RQ1. Correctness**

- **RQ2. Performance**

- **RQ3. Efficiency**

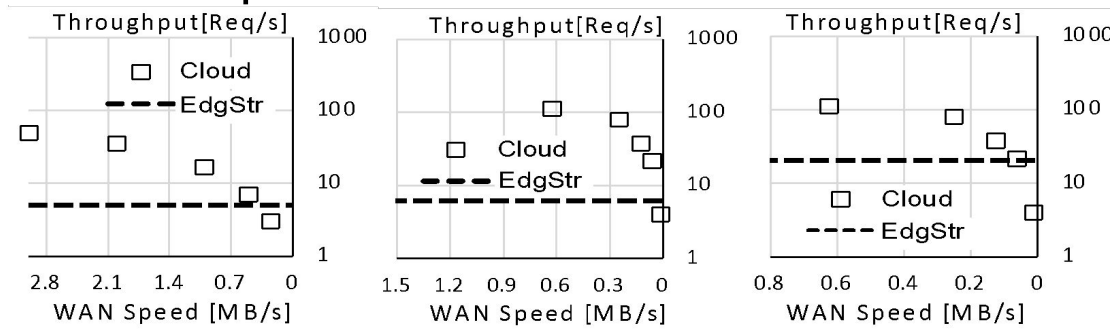  (comparison with related works)

| Components | Specification |
|---|---|
| Cloud Infra (Desktop) | i7-7700 (3.6GHzX8) |
| Edge Node (RPI-3) | Cortex-A53 (1.4GHzX4) |
| Edge Node (RPI-4) | Cortex-A72 (1.5GHzX4) |
| Mobile Dev (Android) | Snapdragon -616 |

# Evaluation (RQ1. Correctness)

- **42/42 was correctly transformed**

    - Given $(p_1, ...p_n)$ sent to the original service OS and the
      replicated service RS, check if Ros == Rrs
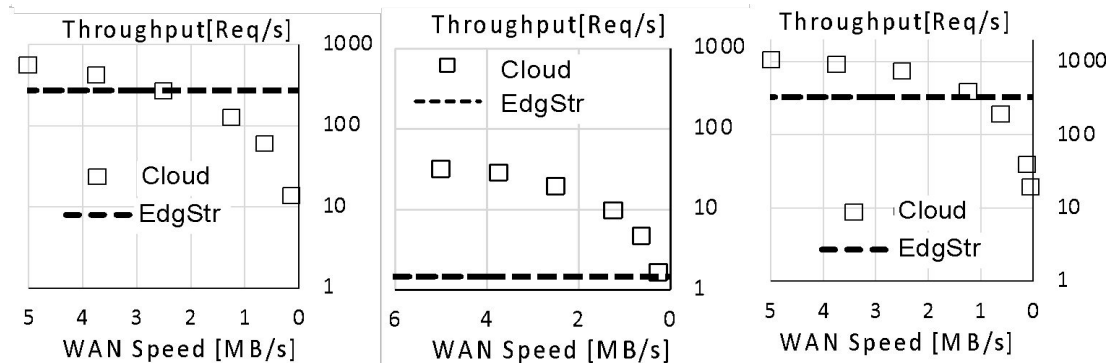
# Evaluation (RQ2: Throughput)

- **Benefit of Edge-based execution** in subjects with
  - Relatively heavy upload/download
  - Low computational loads



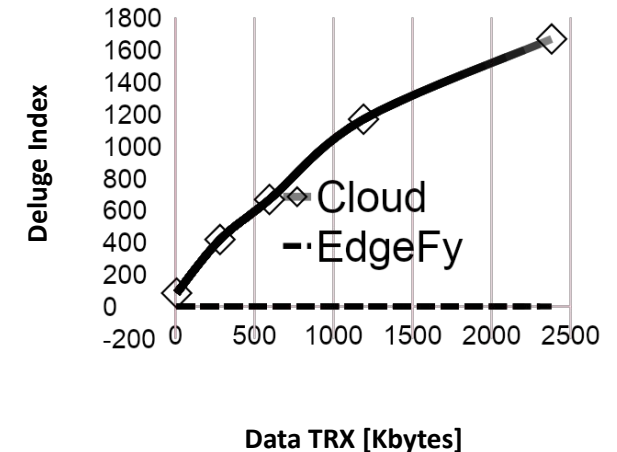(a) f-objdet     (b) mnist-rest     (c) med-chem-rules

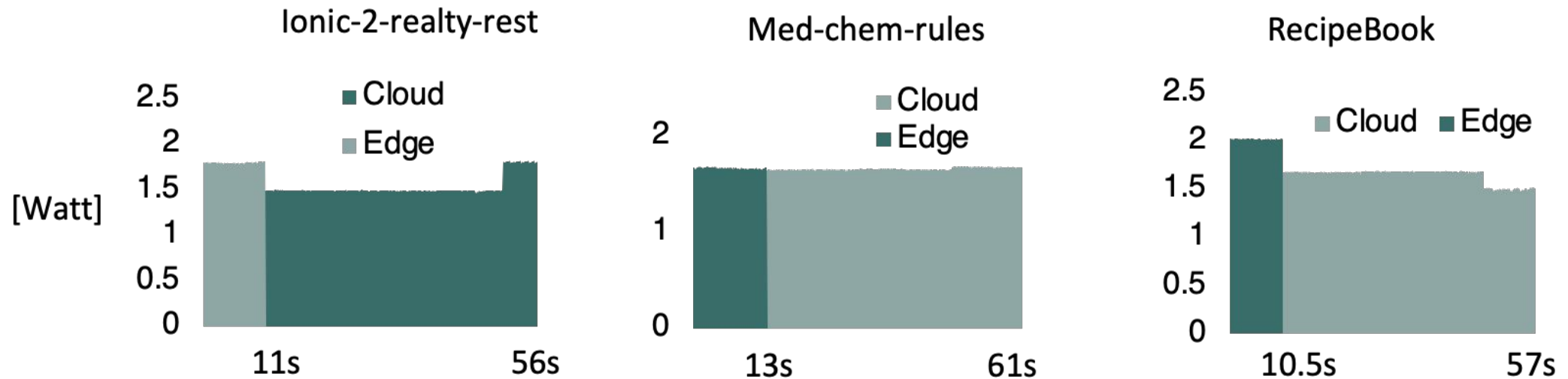(d) ionic2-realty-rest     (e) Bookworm     (f) RecipeBook

- **Deluge Index** ($\Delta$Net/$\Delta$Tput)
  - The volumes of transmitted data over WAN almost did NOT affect EdStr's throughput



23

# Evaluation (RQ2. Energy Consumption in Client Device)

- The longer it takes to execute a cloud-based, the more client device will end up consuming
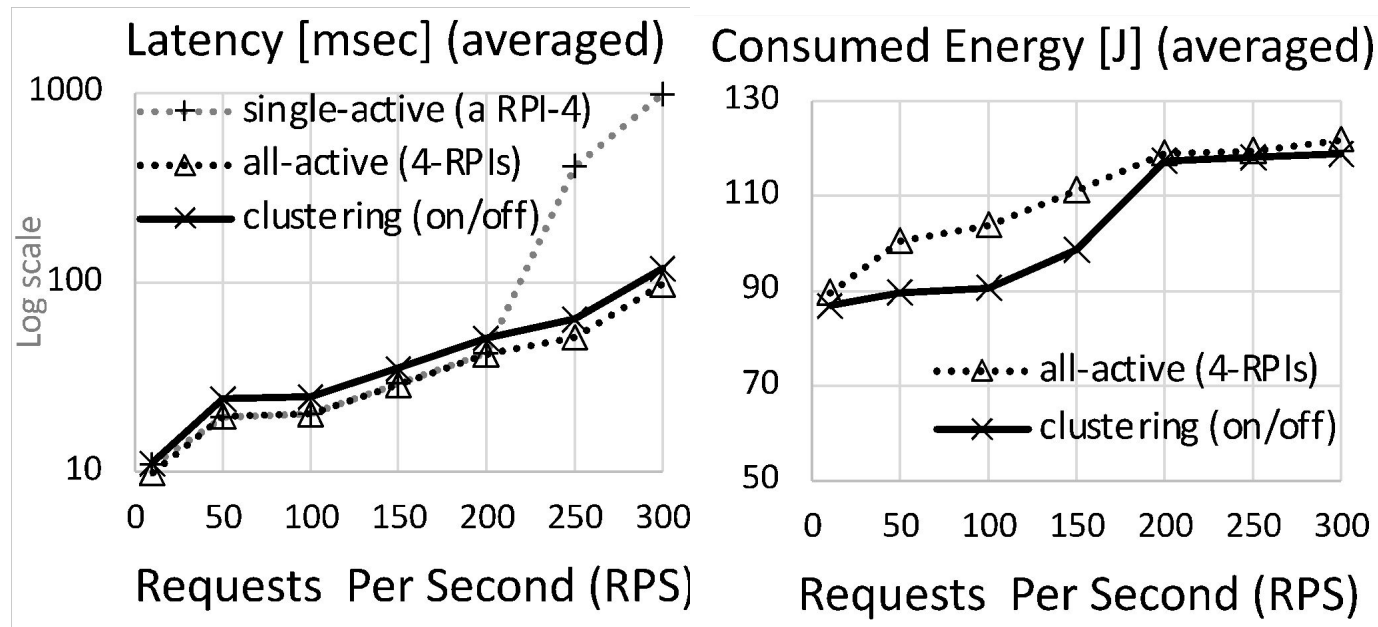


We used Trepn Profiler to measure the consumed energy in Android Device

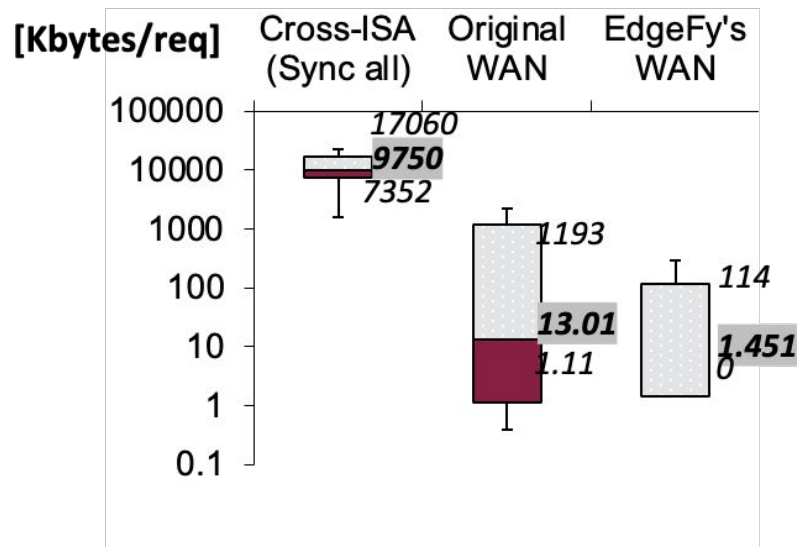# Evaluation (RQ2. Scalability and Elasticity of Edge-based processing)

- Built a cluster using 4 RPIs: distributing clients' requests to available edge replicas
  2 RPI-3s and 2 PRI-4s

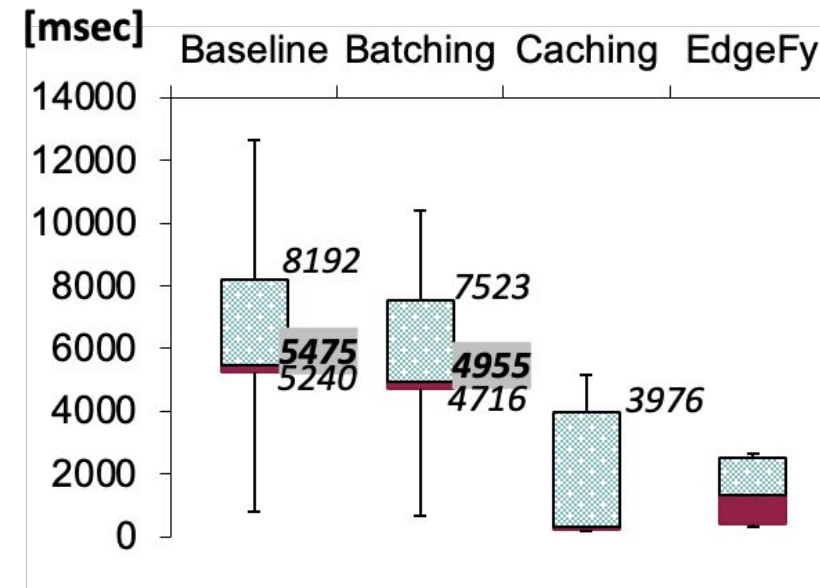- Load balancer shuts on or off the RPIs based on service utilization (clustering on/off)



Active replicas gradually changed from 4 to 1, reducing overall consumed energy by as much as **12.96%**

# Evaluation (RQ3. effectiveness of EdgStr's sync and proxying strategy)

- **Cross-ISA** offloading systems [25,26,27] inefficiently **syncs all states** of cloud program

- Proxy Caching [28,29] benefits in **read-mostly services**

- Batching [31,32] only **reduces WAN traffics** through request aggregation



Sync Overhead and WAN traffic analysis



Comparing the Latency of proxy strategies

Cross-ISA vs EdgStr: EdgStr minimizes the amount of synchronization traffic over WAN by synchronizing only the *modifiable* parts of the replicated service state.

# Conclusion and Q/A

- We described and evaluated EdgStr's advanced program analysis and transformation techniques
  - from 2-tier client-cloud to 3-tier client-edge-cloud

- Applying EdgStr to representative distributed mobile apps introduces the performance benefits of edge processing, without the high costs of manual program transformation

# Applicability & Limitation for **EdgStr**

- Subject: Cloud Services (targeting important domain in Node.js)

- RESTful HTTP protocols

  - Executions: HTTP Request/Response, GET/POST/…

  - What else? Socket.IO, gRPC, …

- ☞ **Cloud Server State Replications**

  - DataBase with **SQL**, **Files**, and **global variables**

  - What else? (Future work) framework specific Data Structures or **ML Models**

    - Federated Learning for replicating **ML Models** across cloud and edges