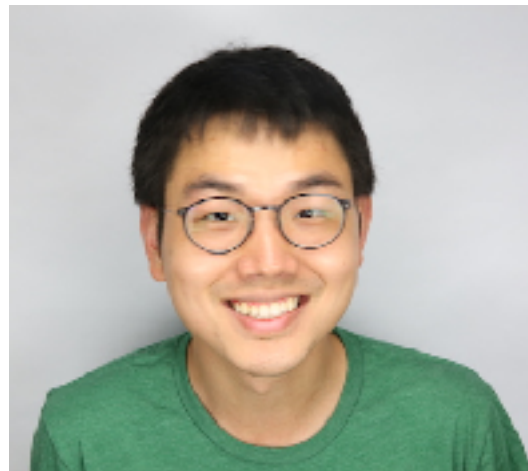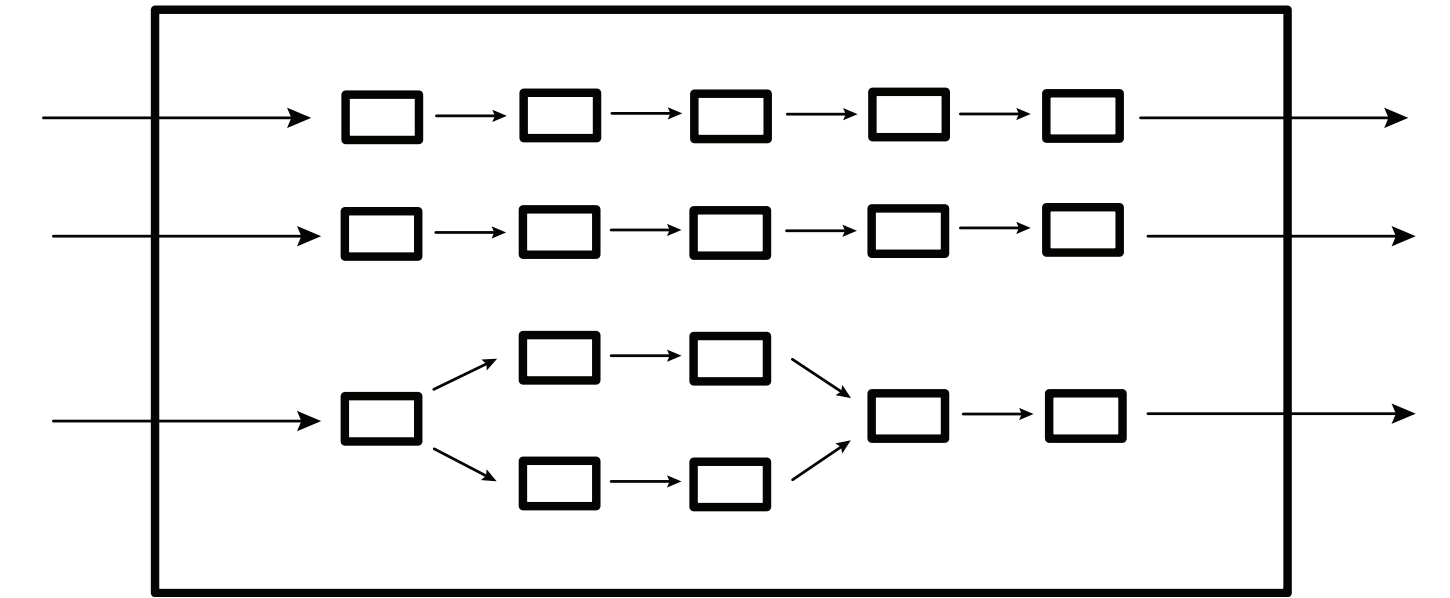# *Peekaboo*

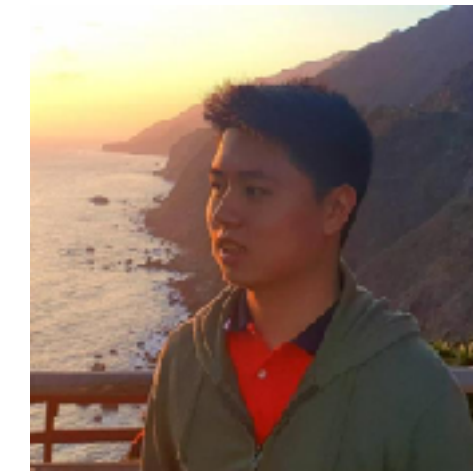## A Hub-Based Approach to Enable Transparency in Data Processing within Smart Homes



**Haojian Jin**

Gram Liu

David Hwang

Swarun Kumar
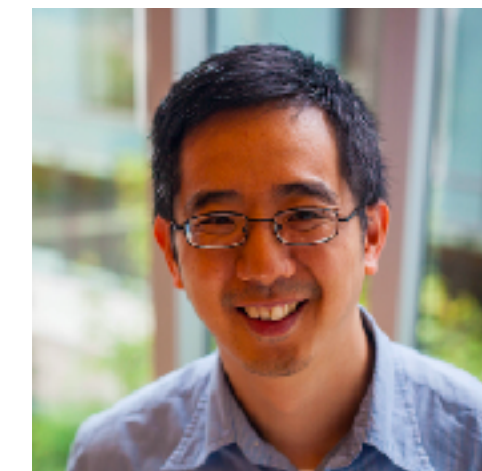
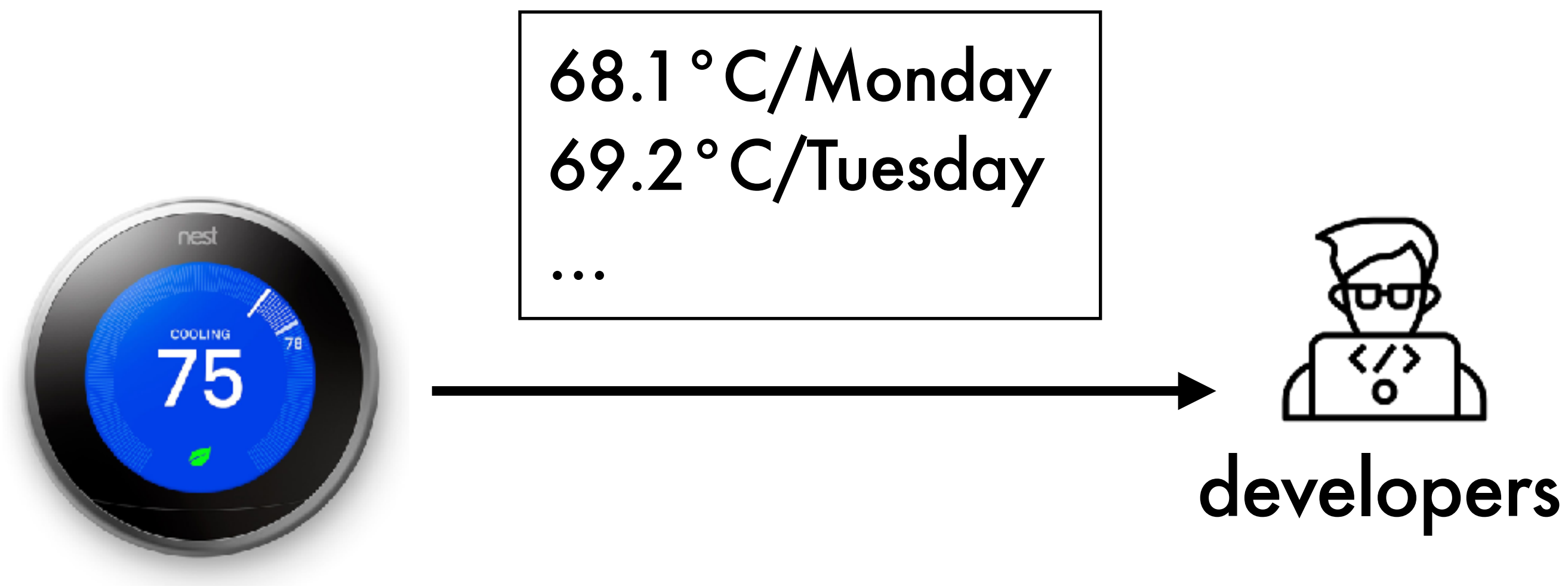Yuvraj Agarwal

Jason Hong

# How can Nest prove that they only collect aggregated data?



68.1°C/Monday
69.2°C/Tuesday
...

developers

Open source?

# Your TV watch history contains too much insights

| video # | duration | name | time | ... |
|---------|----------|------|------|-----|
| aaa | - | - | - | - |
| bbb | - | - | - | - |
| ... | ... | ... | ... | ... |

25 hours/week

How much time does the user spend on the TV?

- Is the user at home
- Activity routine
- User interests
- ....

# Only collect the necessary data for a specific purpose.
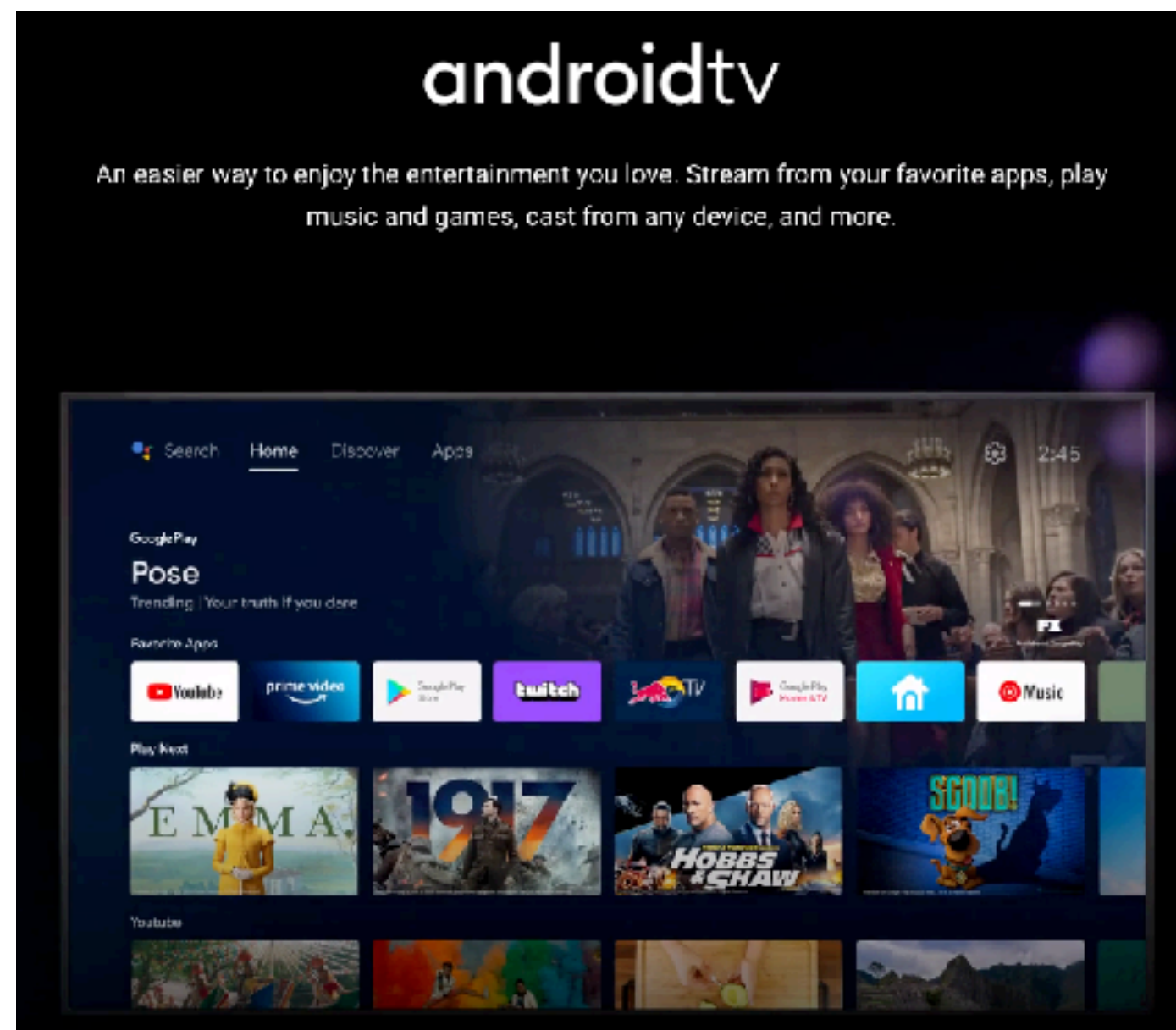


25 hours/week

Use weekly usage data to measure device engagement.

- ~~Is the user at home~~
- ~~Activity routine~~
- ~~User interests~~
- ....

How can developers prove themselves?

4

# A strawman solution: fine-grained permission manifest



https://www.android.com/tv/

```
<manifest ...>
  <uses-permission android:name="android.permission.
      TV_AGGREGATED_DURATION_WEEKLEY" />

<uses-permission android:name="android.permission.
      TV_AGGREGATED_DURATION_DAILY" />


  ......
</manifest>
```
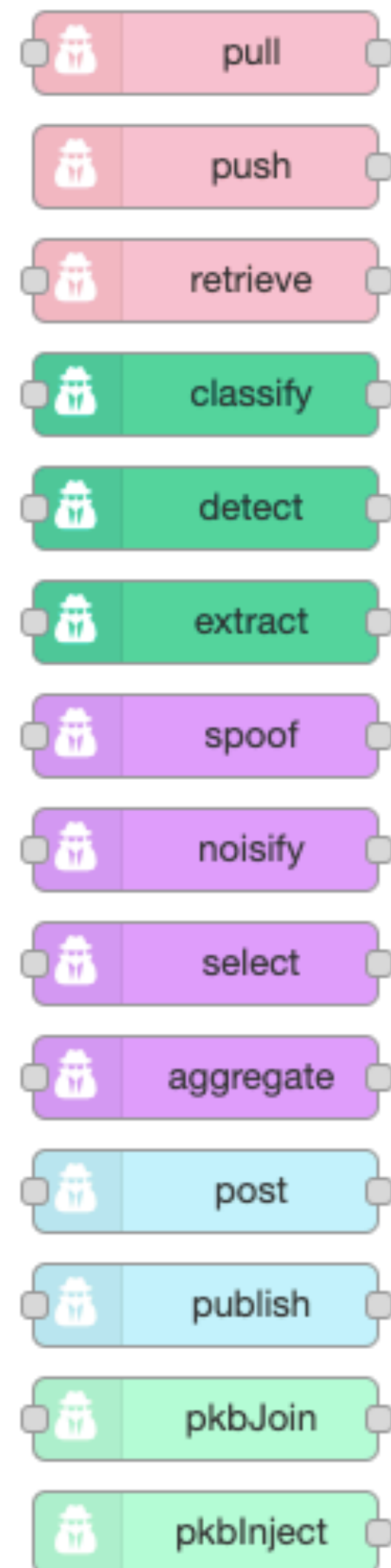
Fine-grained permission manifest

# Program pre-processing functions using chainable *operators*



A fixed set of operators

pull
push
retrieve
classify
detect
extract
spoof
noisify
select
aggregate
post
publish
pkbJoin
pkbInject

inject [weekly] — pull smart tv — aggregate [sum duration] — post [duration abc.com]

Edit aggregate node

Delete          Cancel    Done

⚙ Properties

🏷 Name          aggregate [sum duration]

⛁ Data Type      tabular

⊙ Target         custom
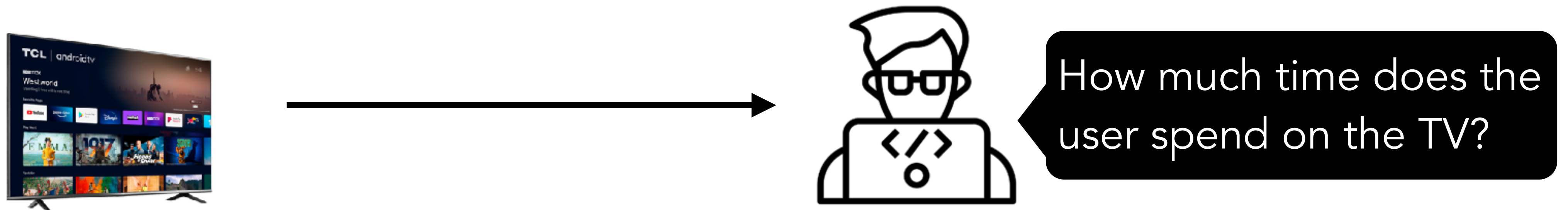
▤ Tabular field  duration

▤ Operation      sum

⚙ Options
(optional) ▾

⚙ Group by

▤ As ?           new variable name

6

# A text-based whitelist *manifest* (i.e., program representation)



How much time does the user spend on the TV?

```
@purpose: To measure device engagement.
WeeklyUsageHours{
    // operator [properties]
    inject [weekly] ->
    pull [smart TV driver] ->
    aggregate [sum duration] ->
    post [duration]
}
```
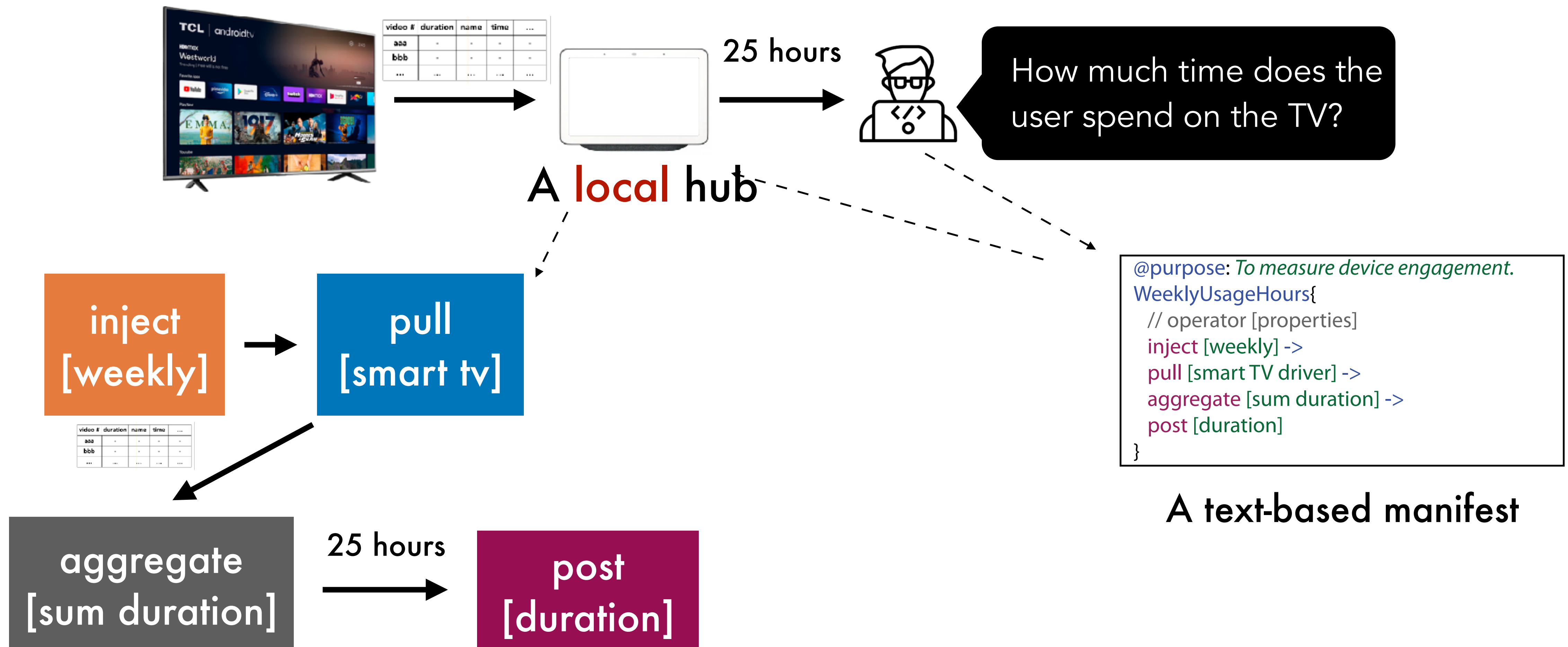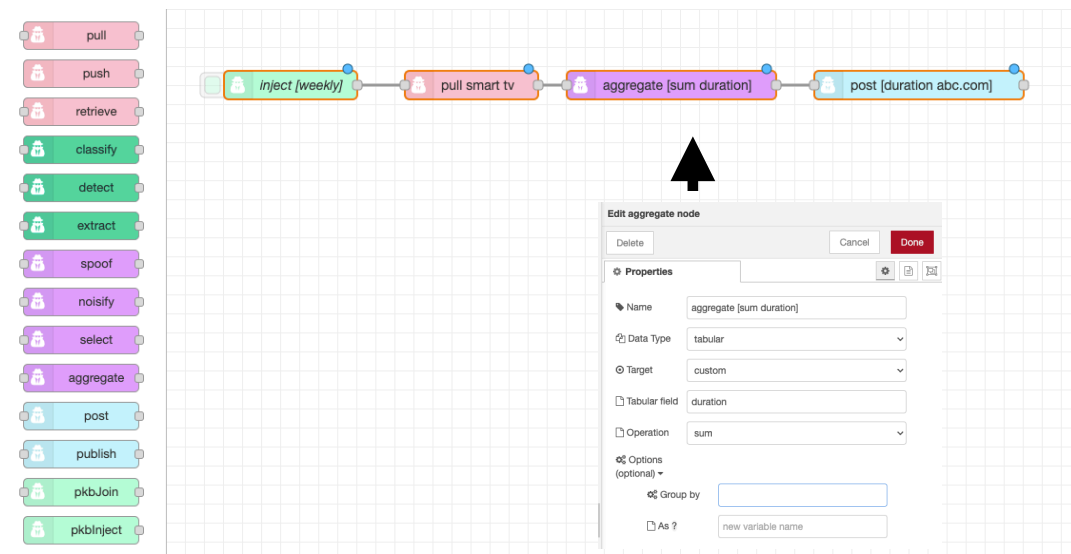
# A trusted runtime with pre-loaded implementations



25 hours

How much time does the user spend on the TV?

A local hub

inject [weekly]

pull [smart tv]

aggregate [sum duration]

25 hours

post [duration]

```
@purpose: To measure device engagement.
WeeklyUsageHours{
  // operator [properties]
  inject [weekly] ->
  pull [smart TV driver] ->
  aggregate [sum duration] ->
  post [duration]
}
```

A text-based manifest

8

# Smart home app store

# App developers



## Programming environment with operators

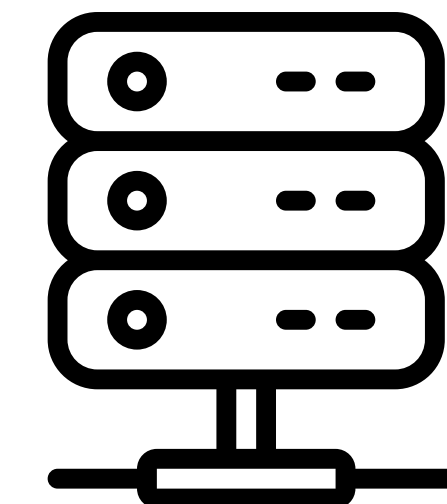## Runtime with preloaded implementations

```
@purpose: To measure device engagement.
WeeklyUsageHours{
  // operator [properties]
  inject [weekly] ->
  pull [smart TV driver] ->
  aggregate [sum duration] ->
  post [duration]
}
```
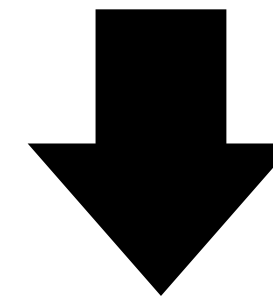
## Manifest

# Smart home app store

```
@purpose: To measure device engagement.
WeeklyUsageHours{
  // operator [properties]
  inject [weekly] ->
  pull [smart TV driver] ->
  aggregate [sum duration] ->
  post [duration]
}
```

Smart home app →
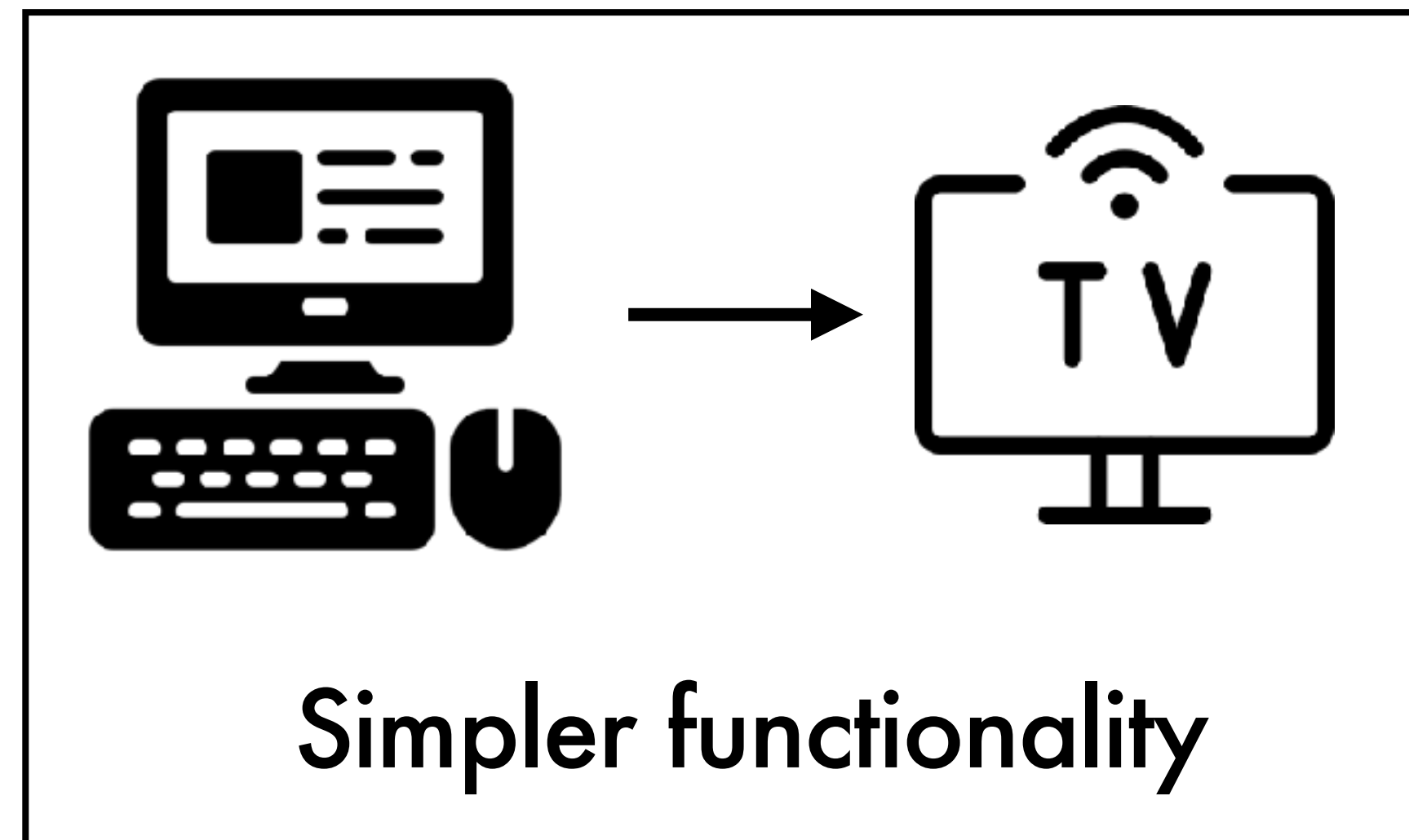
Edge devices

A local hub

Cloud

*"Privacy firewall"*

10

# Developers declare purposes explicitly.

Simpler functionality

@purpose: *To measure device engagement.*
WeeklyUsageHours{
    p        p   p

**Whitelist-only**

**Developer-in-the-loop**

# 77% Apps do not need raw data.

| | Sensor | Raw | Needed data |
|---|---|---|---|
| Hello visitor | | | |
| Noise level | | | 55 db |

# Pre-process users' data to mitigate data overaccess.

| video # | duration | name | time | … |
|---------|----------|------|------|---|
| aaa | - | - | - | - |
| bbb | - | - | - | - |
| … | … | … | … | … |

25 hours/week

Edge devices

Hub

Cloud

# Recap: Peekaboo v.s. Firewall

Simpler functionality

➡️ Whitelist-only

Developer-in-the-loop

77% Apps do not need raw data.

➡️ Pre-process users' data

# Handle heterogeneous hardware with device drivers



Device APIs

Edge devices

Device drivers

inject

pull

aggregate

post

HomeAssistant, https://www.home-assistant.io/

# A fixed set of operators



Edge devices

video, image, audio, tabular, scalar

A *fixed* set of operators

pull
push
retrieve
classify
detect
extract
spoof
noisify
select
aggregate
post
publish
pkbJoin
pkbInject

# An operator = A verb keyword



**select [row]**

| | product_id | product_name | inventory_received | starting_inventory | inventory_on_hand | minimum_required |
|---|---|---|---|---|---|---|
| 1 | 2 | Booth | 29pcs | 27pcs | 56pcs | 20pcs |
| 2 | 3 | Maclean | 23pkts | 25pkts | 48pkts | 25pkts |
| 3 | 4 | Closeup | 24pkts | 25pkts | 49pkts | 25pkts |

**detect [face]** → **select [face]** →

# Operators are mapped to pre-loaded implementations



select
[row]

select
[face]

Row selection

Image cropping

# A small set of pre-processing algorithms improve privacy

| video # | duration | name | time | ... |
|---------|----------|------|------|-----|
| aaa | - | - | - | - |
| bbb | - | - | - | - |
| ... | ... | ... | ... | ... |

Row selection

Image cropping

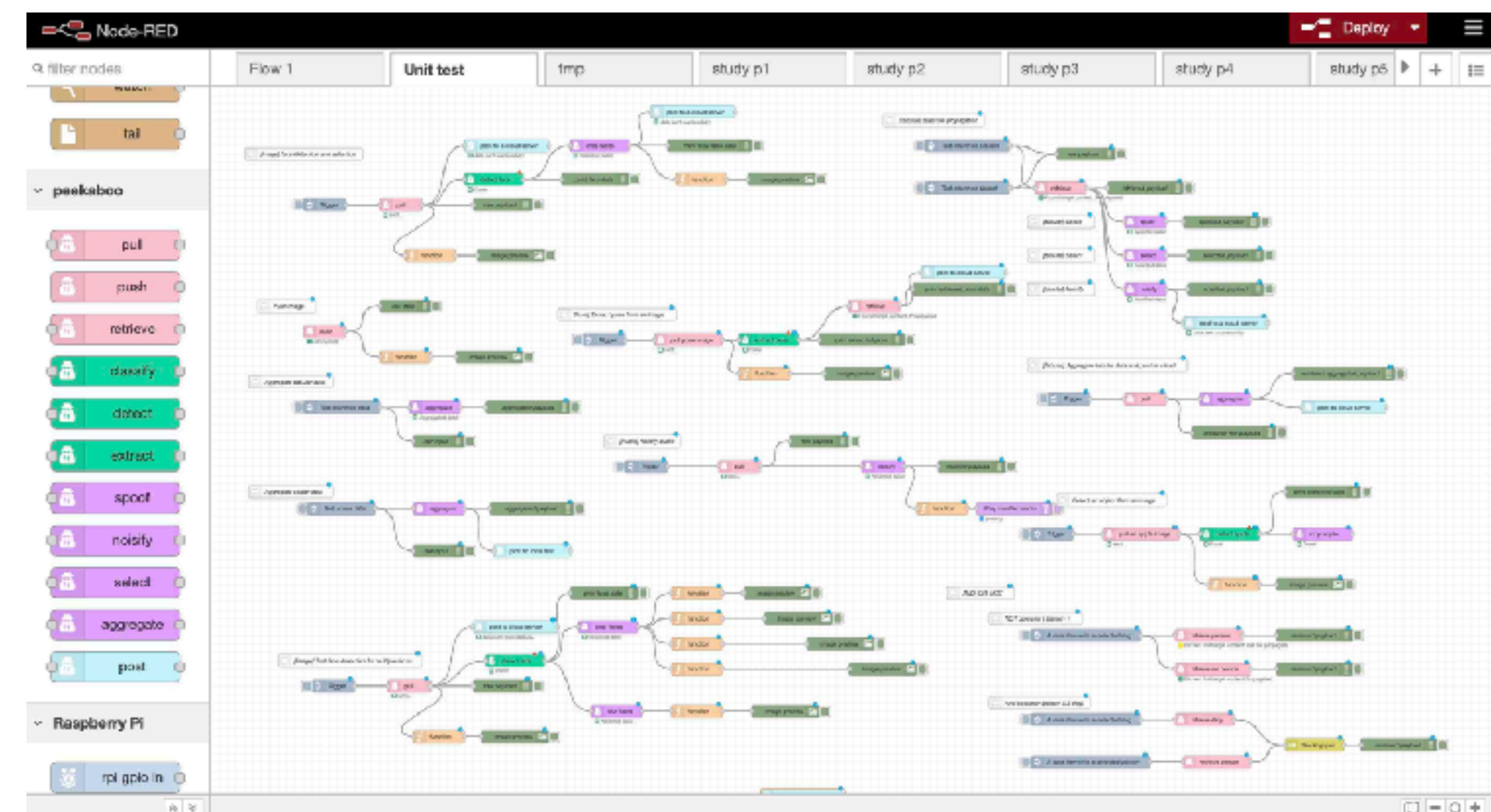25 hours/week

# Implementation (hardware)
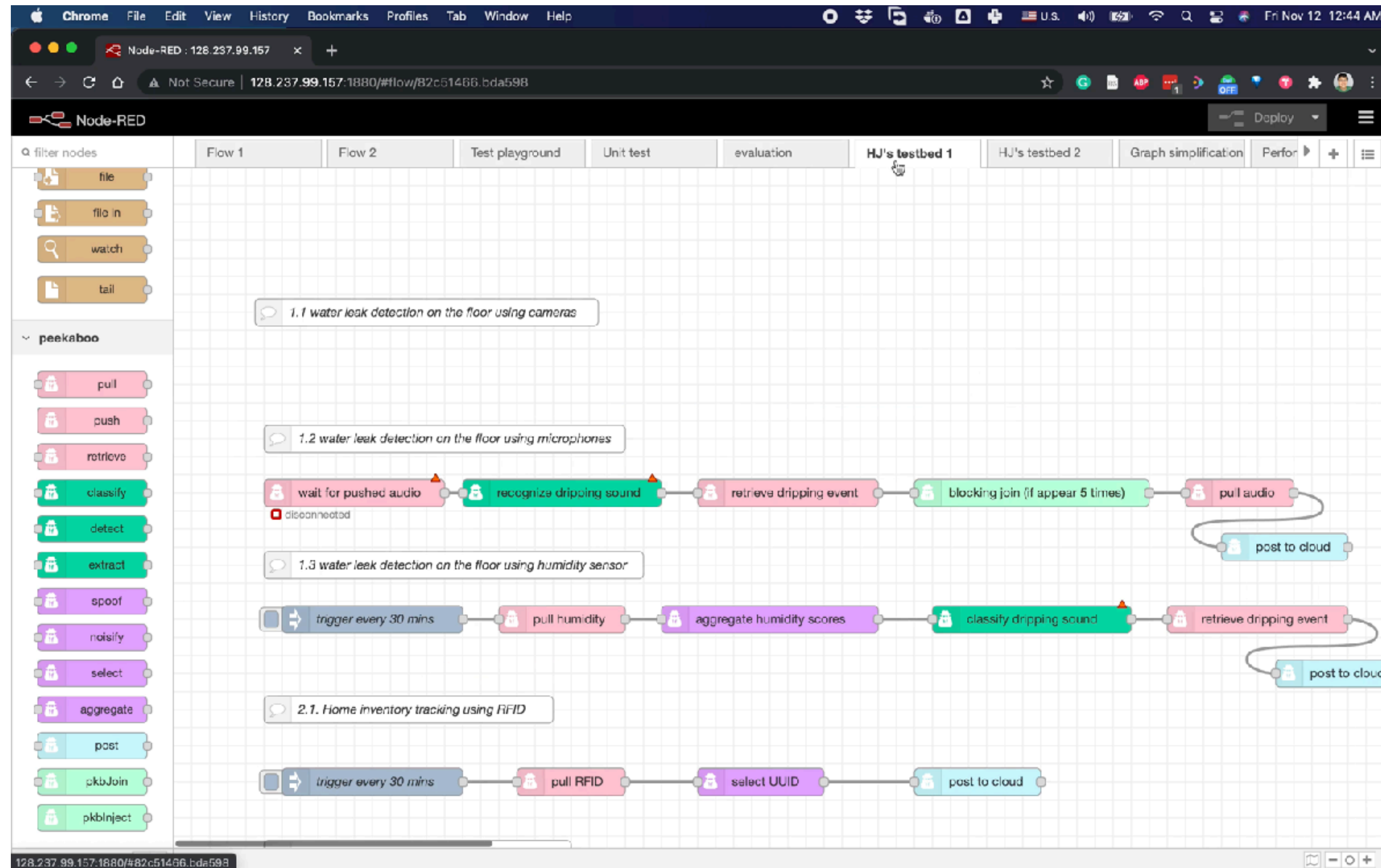

Edge devices


Raspberry PI + TPU


Cloud

# Implementation (software)

1. Operators: Node.JS package

2. Programming IDE: NodeRed

3. Drivers: 5 data types

4. 23 Preloaded implementations

# Expressiveness (200+ smart home cases)

# Data overaccess mitigation breakdown

unique manifests: 68

content selection: 64

explicit noisification: 57

conditional filtering: 51

See details in the paper

**3** cannot mitigate → push → post

# System performance



≈$100

25 inference/s

100 filtering/s

1-80 ms per request

24

# Utility privacy tradeoff example



incognito voice assistant

6 speakers
112 audio files [1]
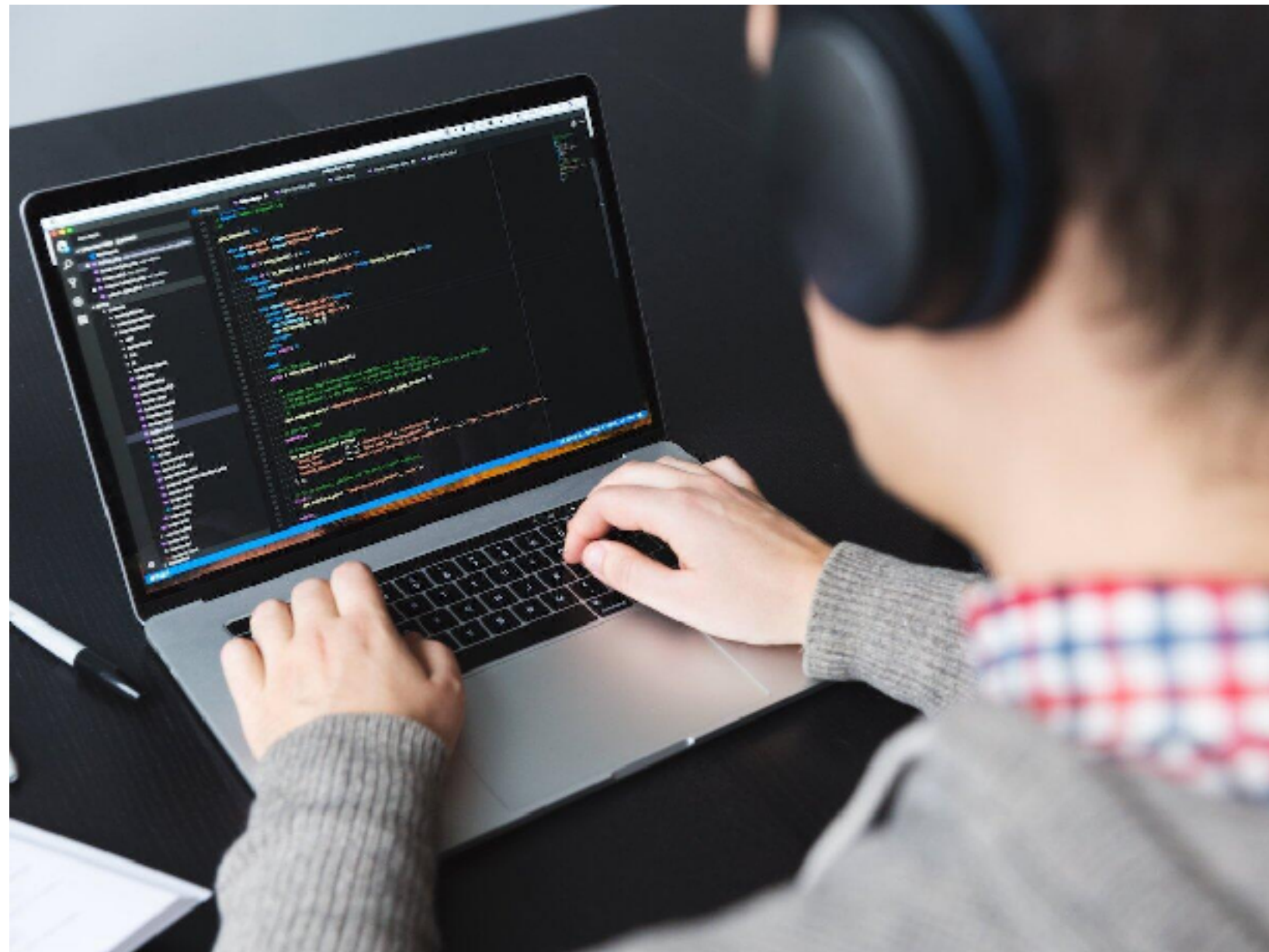
| noisify | <5% random pitch shift |

Speech word error rate:
*9.27% → 11.88%*

Speaker recognization:
*100% → 27.7%*

Microsoft
Cognitive Services

25

[1] CMU PDA Speech Database
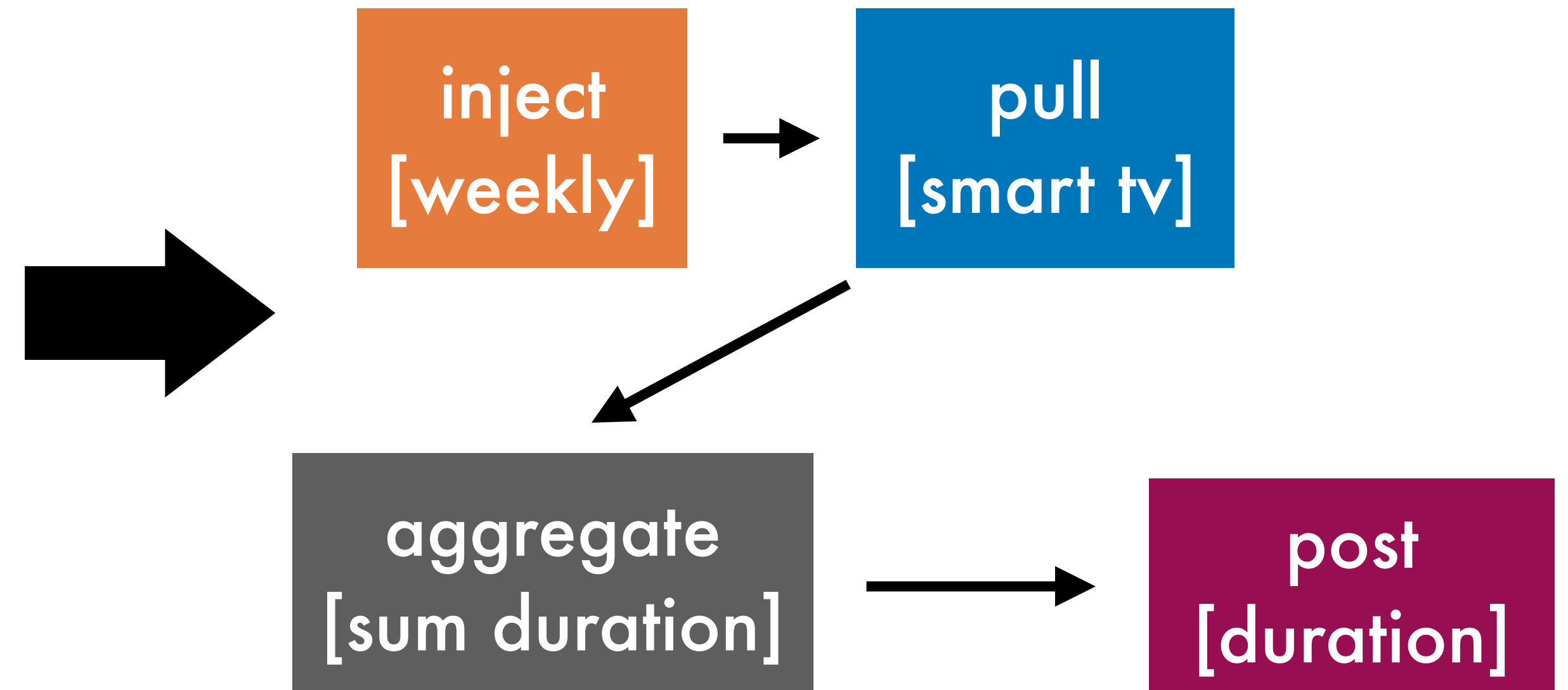
# Developer studies



Task descriptions

IDE & Unit tests

*6 - 15 mins* to
author a manifest

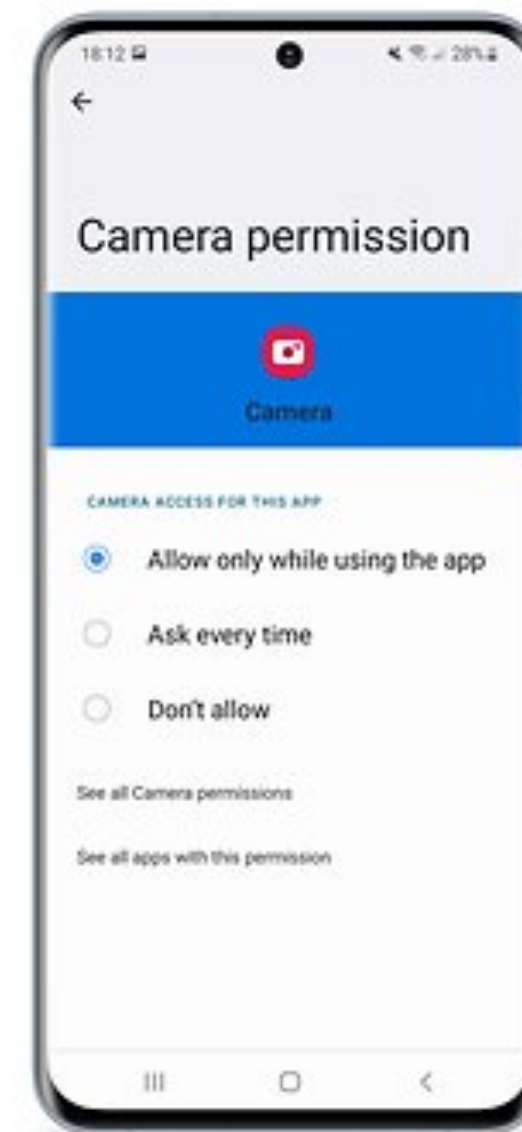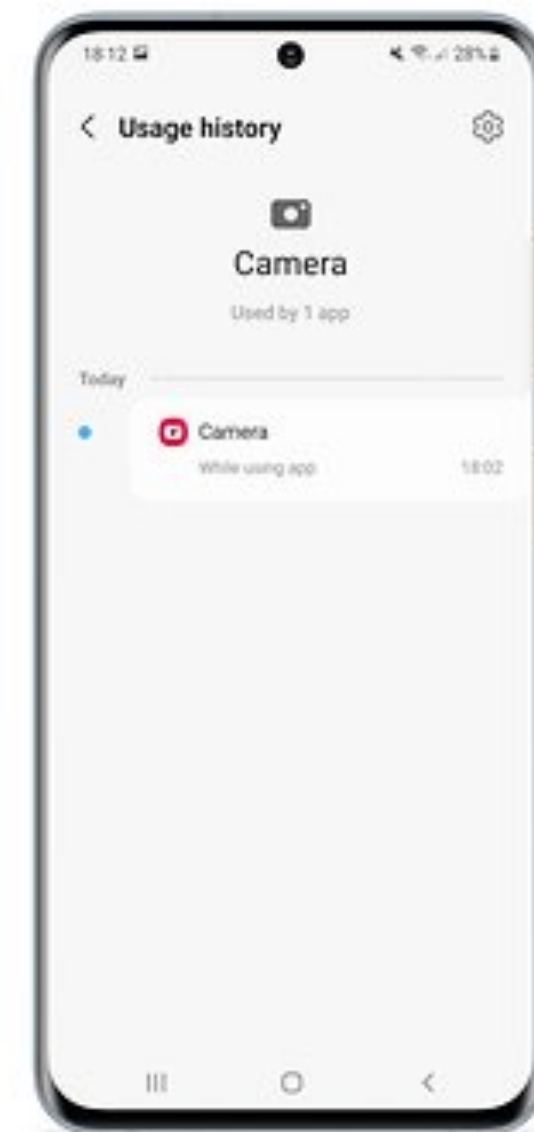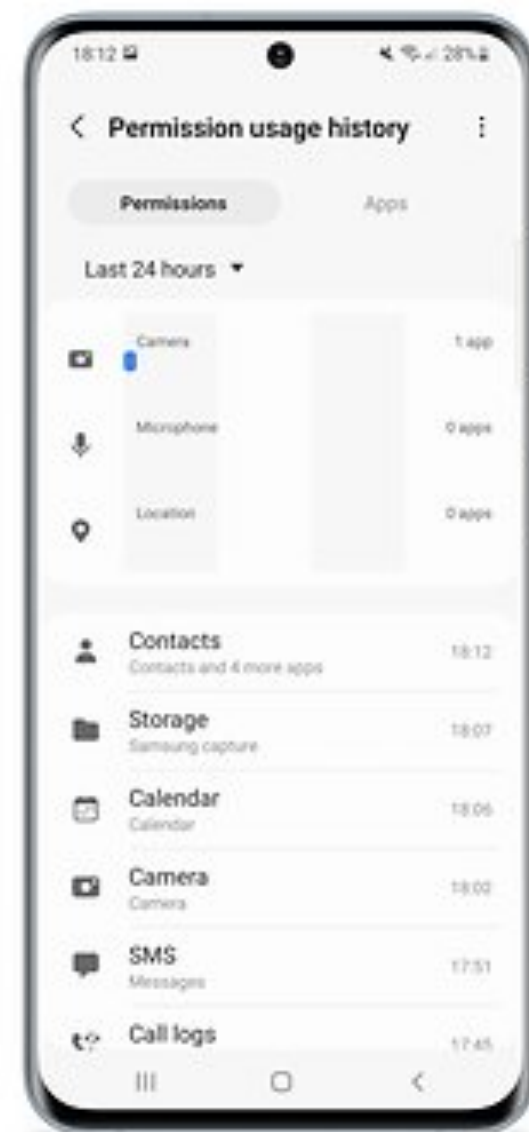# Manifests enforce fine-grained data collection

```
@purpose: To measure device engagement.
WeeklyUsageHours{
  // operator [properties]
  inject [weekly] ->
  pull [smart TV driver] ->
  aggregate [sum duration] ->
  post [duration]
}
```
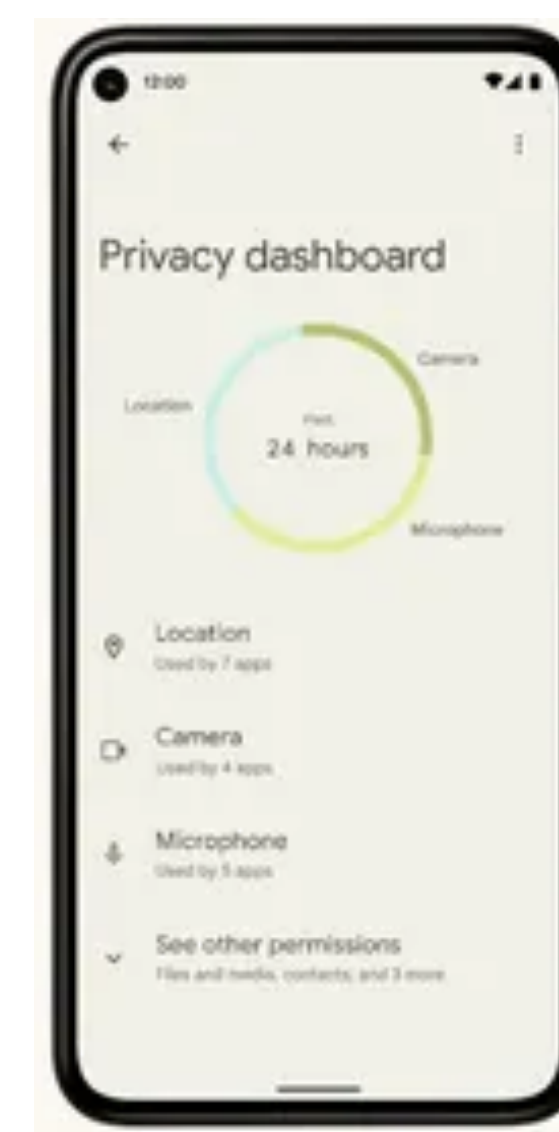
public, non-proprietary



27

# *Repetitive* implementation and *distributed* interfaces

Samsung

Nest



*Small developers?*

*Users?*

# Manifests → *enforceable/dynamic* privacy nutrition labels

[1]

```
@purpose: To measure device engagement.
WeeklyUsageHours{
  // operator [properties]
  inject [weekly] ->
  pull [smart TV driver] ->
  aggregate [sum duration] ->
  post [duration]
}
```

**Data Collection Disclosure**

| | |
|---|---|
| TV Usage Summary App | |
| **Running for** | 20 days |

| | |
|---|---|
| Total outgoing data packets | |
| **KBytes** | **80** |

| | |
|---|---|
| Sensor Type | Smart TV |
| Data type | TV Watch history |
| Granularity | Weekly aggregated durations by content category |
| Collection frequency | Every wednesday 1:00 AM |
| Destination | www.abc.com |
| Encryption | HTTPS |

**Customizations**

| | |
|---|---|
| Rate limiting | N/A |
| More options | .... |

29

[1] Security and Privacy "Nutrition" Label, P. Emami-Naeini et al, IEEE S&P'20
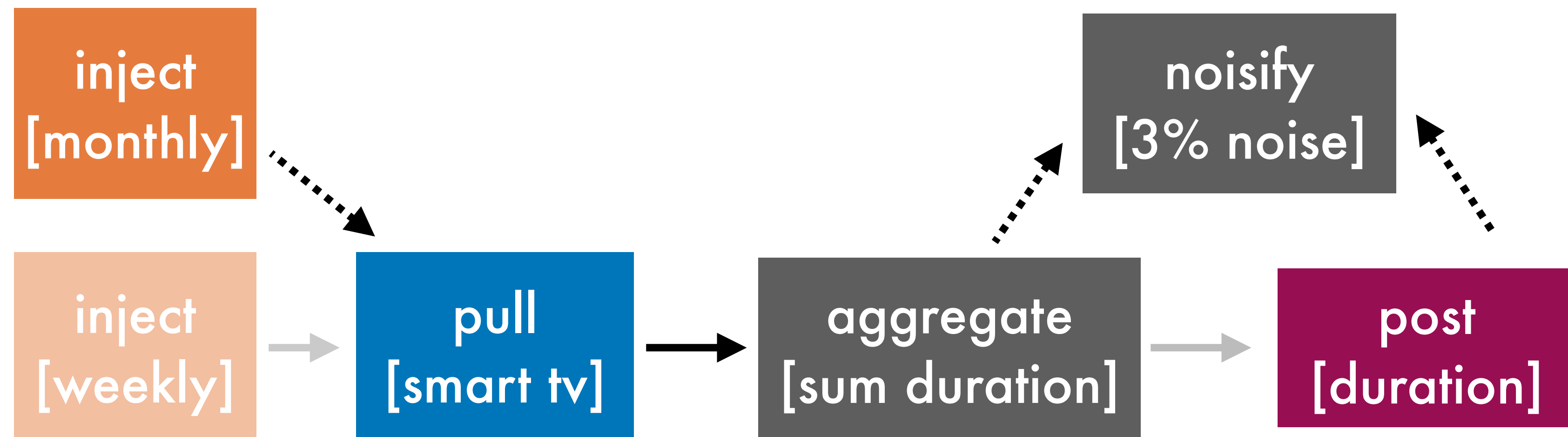
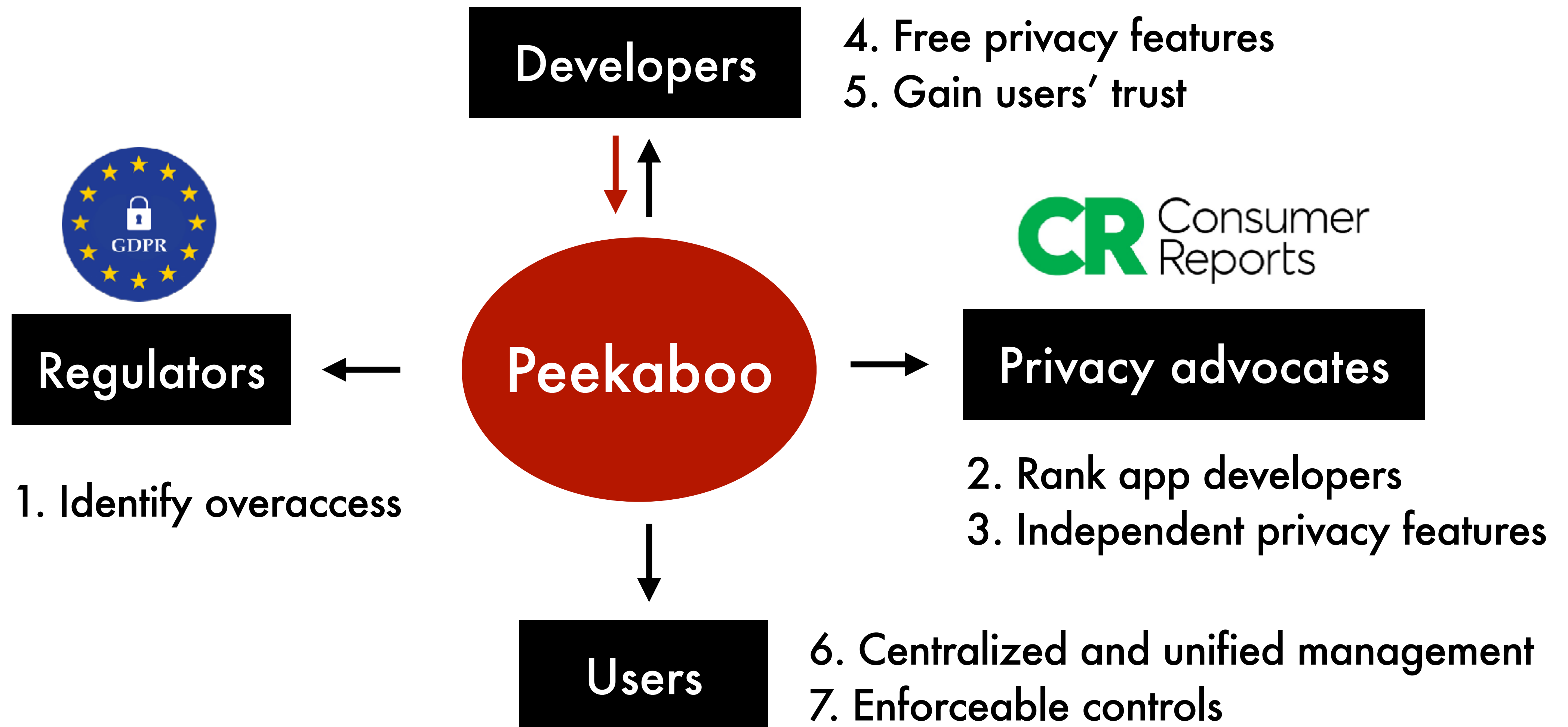# *Built-in fine-grained* control through manifest rewriting



**Data Collection Disclosure**

TV Usage Summary App

| | Customizations |
|---|---|
| Rate limiting | N/A |
| More options | .... |

Change the rate
to monthly

inject
[monthly]

inject
[weekly]

pull
[smart tv]

aggregate
[sum duration]

noisify
[3% noise]

post
[duration]

# Let the good privacy drive out the bad privacy



Developers

4. Free privacy features
5. Gain users' trust

Regulators

Peekaboo

Privacy advocates

1. Identify overaccess

2. Rank app developers
3. Independent privacy features

Users

6. Centralized and unified management
7. Enforceable controls

# Design data access for third-party developers

URL-based APIs



Operator-based APIs

inject → pull → detect → select → post

# Peekaboo recap & implications

Manifest

↓

Operators

A fixed set of operators

↓

Runtime

A trusted runtime with a small set of pre-loaded implementations

↓

Executable

# Zoom accesses **all** your calendar events **continuously**!



Calendar events that contain
https://zoom.us/xxxxx

34

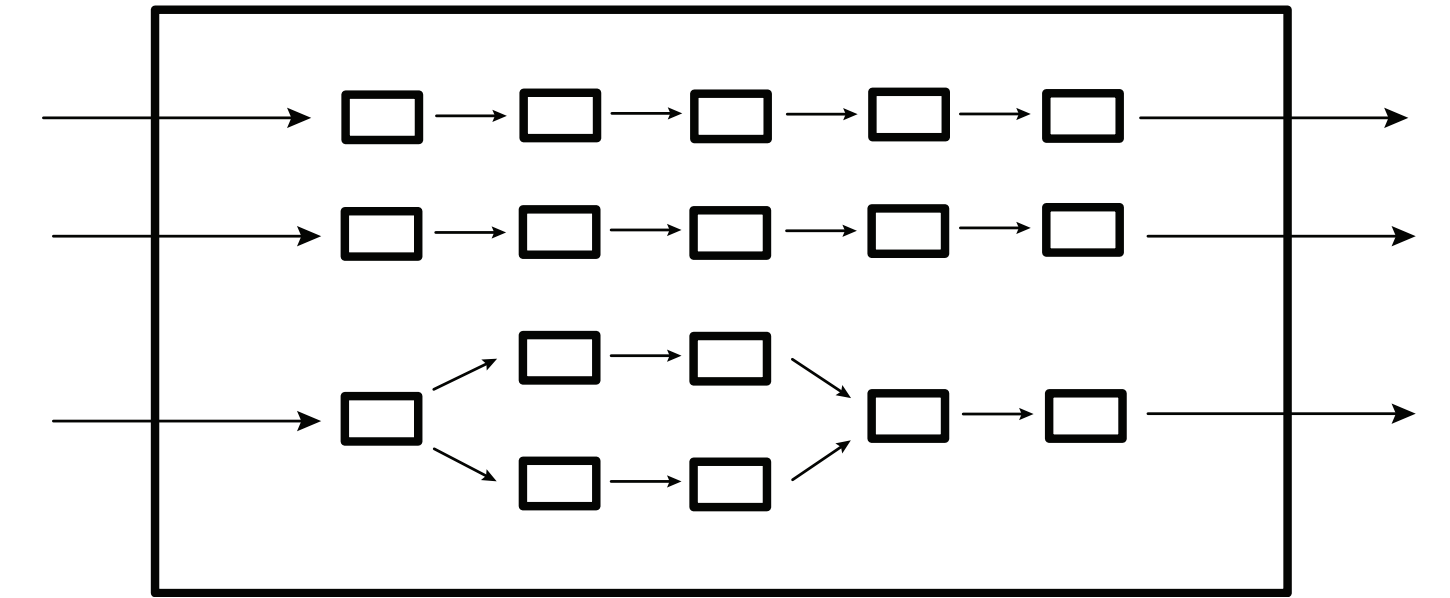# Future third-party calendar API



inject | pull | select | check

...

@purpose: *The app can access calendar events which contains a zoom link.*
ZoomCalendarIntegration{
  // operator [properties]
  inject[...] -> pull Calendar[...] ->
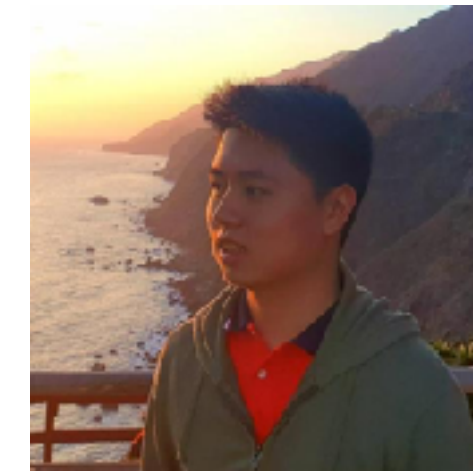  check Zoomlink[...] ->
  post [Zoom events]
}

# Principle of data minimization

"*Personal data shall be limited to* **what is necessary** *in relation to the* **purposes** *for which they are processed.*"

- GDPR, Article 5 (1) (c)