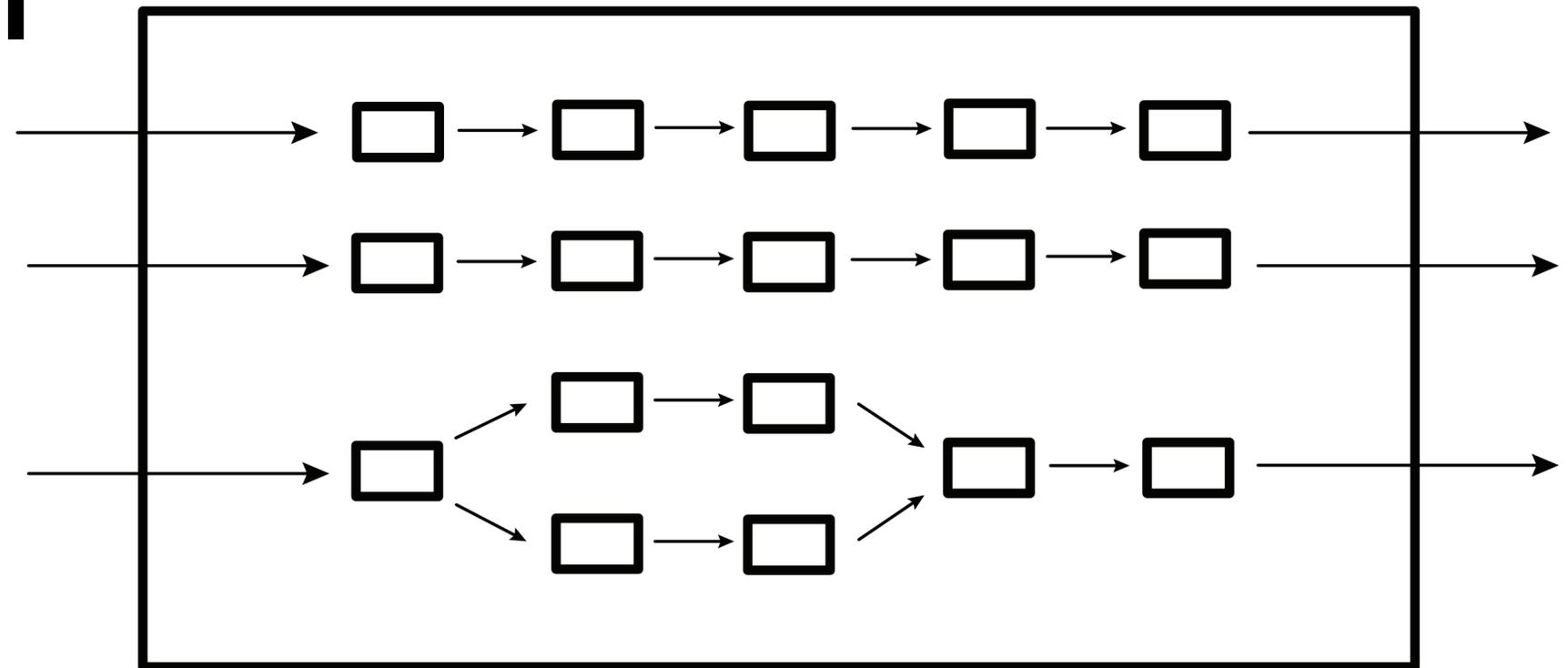


Modular Privacy Flows: A Design Pattern for Data Minimization



Haojian Jin
Mar. 29, 2023

Bio

Haojian Jin (<http://haojianj.in/>)

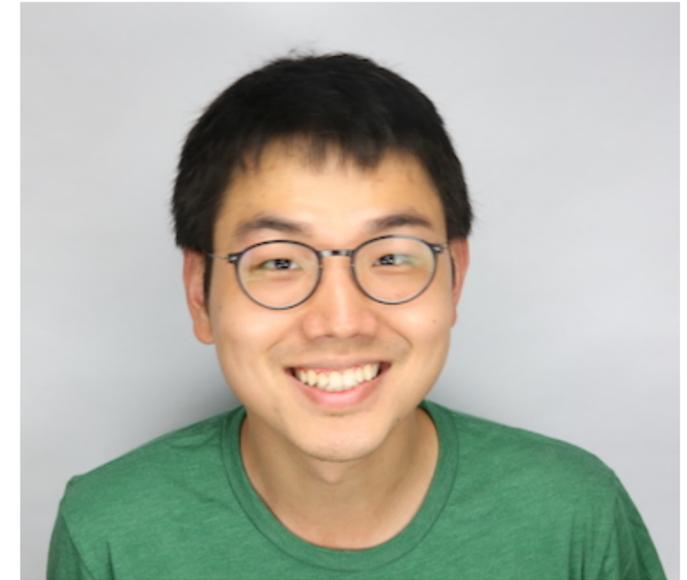
Asst. Prof @ UCSD-HDSI

Data Smith Lab:

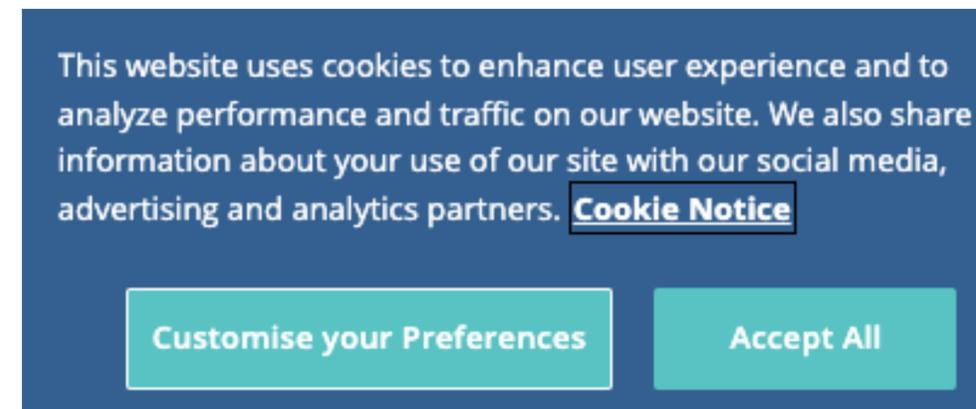
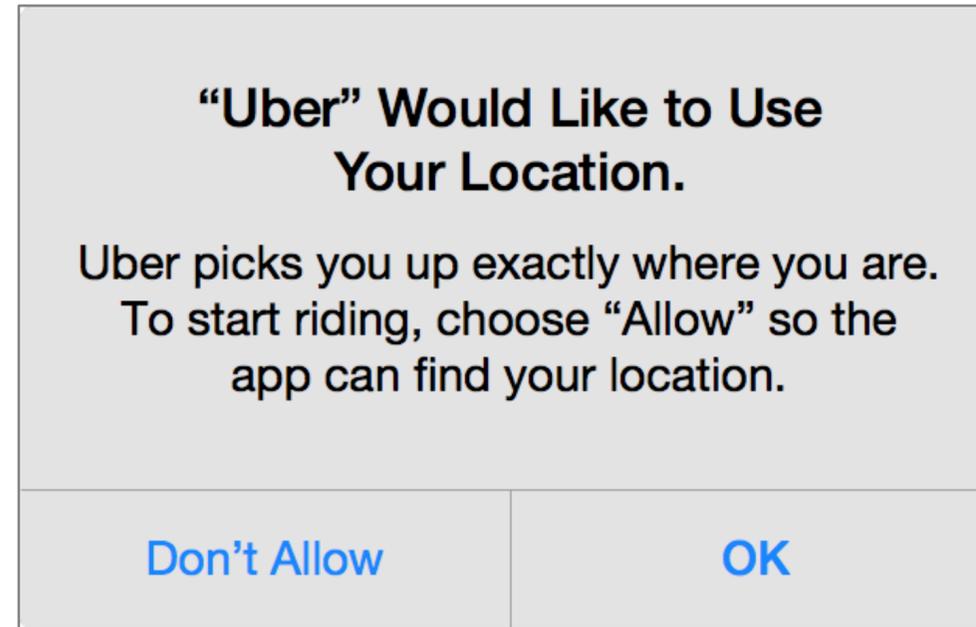
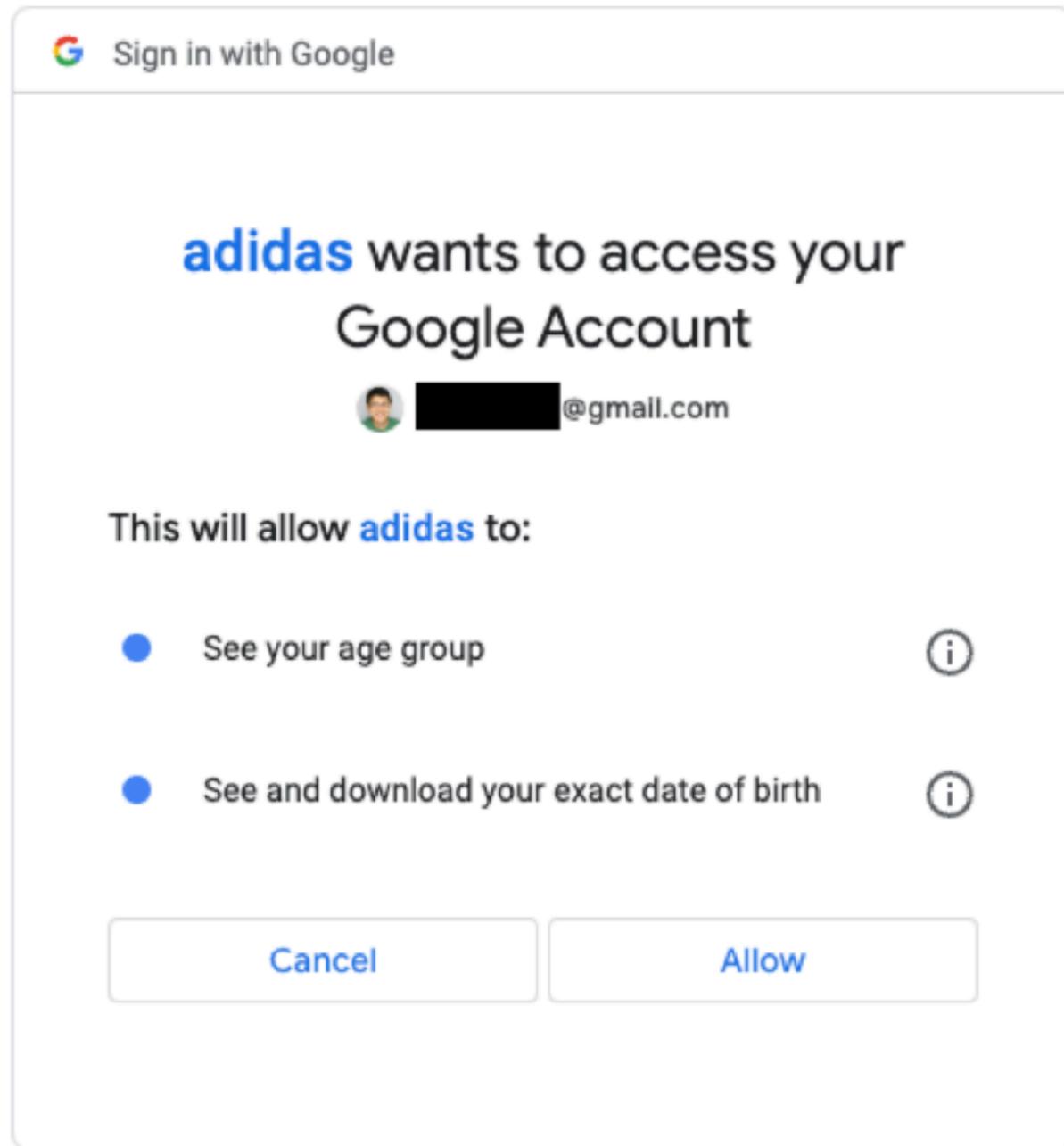
We study the **security and privacy of data systems** by researching the people who **design, implement, and use** these systems.

Ph.D. from CMU Human-Computer Interaction Institute

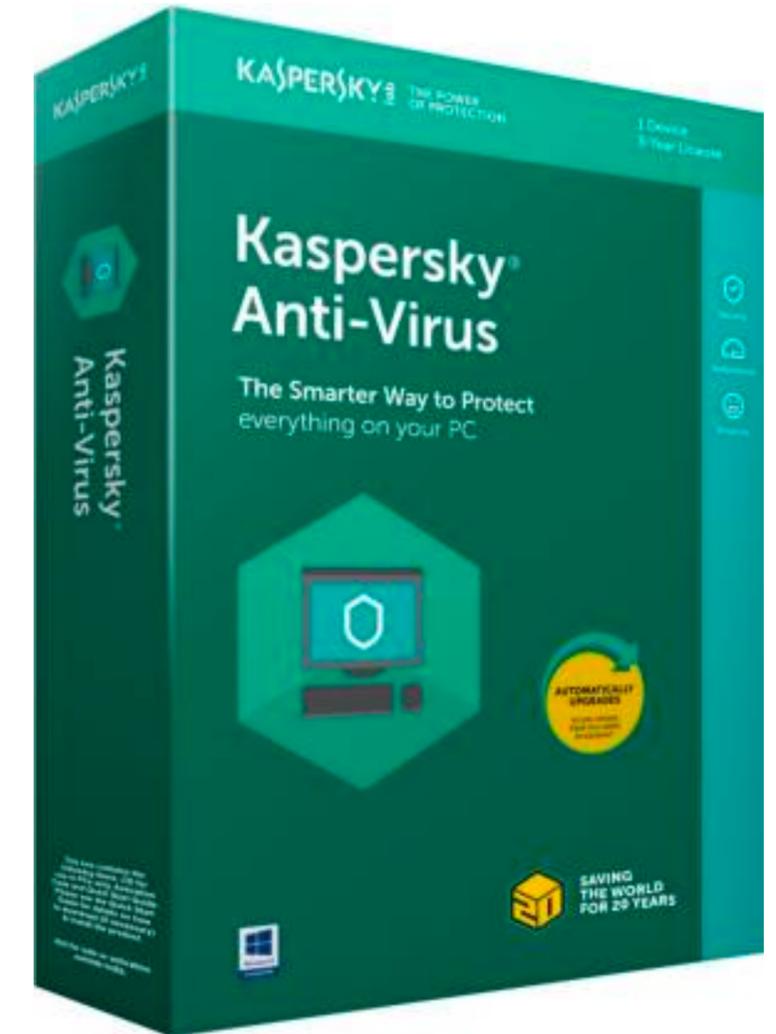
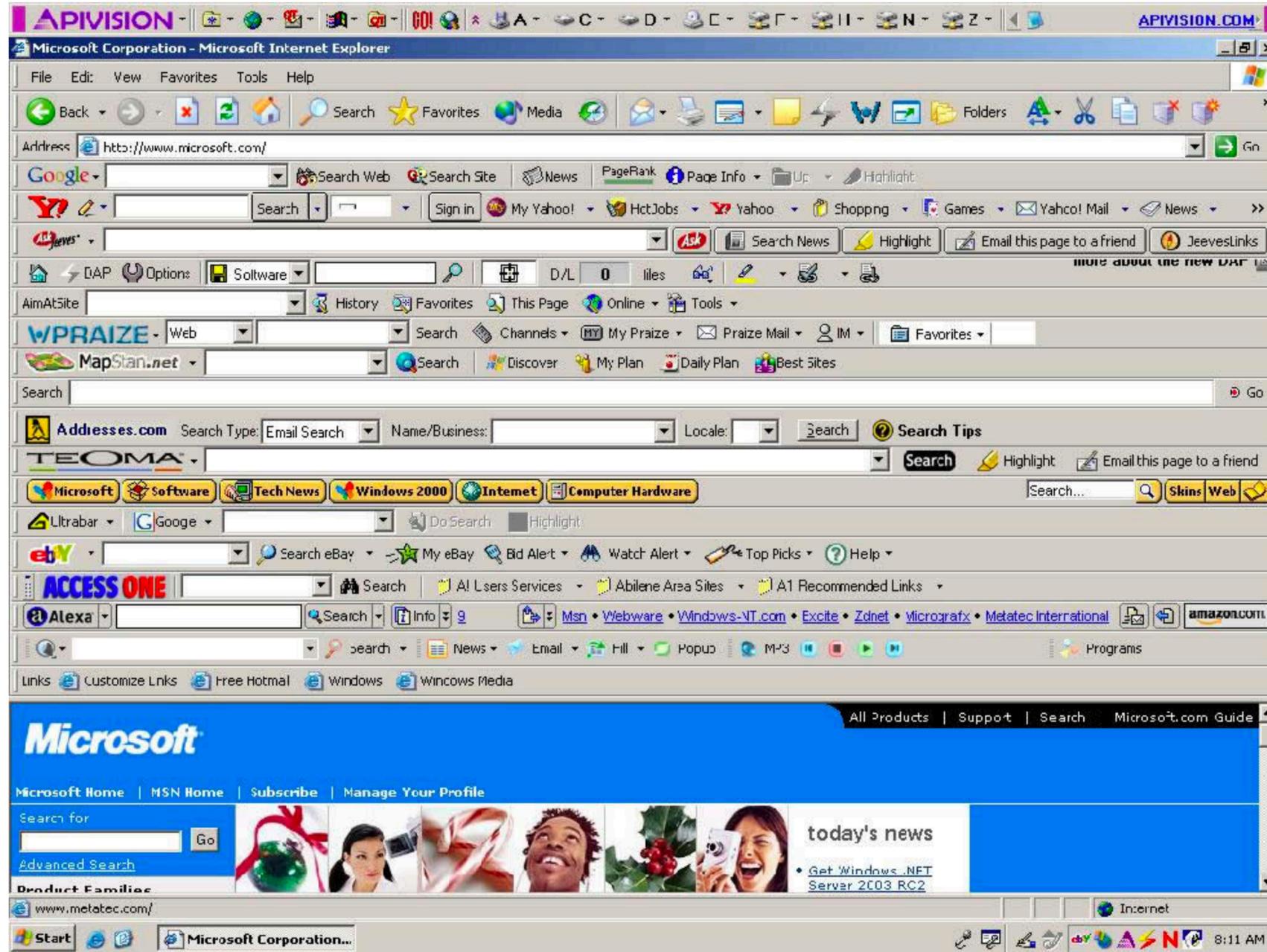
Before Ph.D.: worked at Yahoo Research, ran a startup



Permissions



Notice and choice
Informed decisions



App Stores

*“changed how software development worked, and expanded the **number of people who could comfortably, safely use a computer** from a few hundred million to a few billion.”*

Technical idea #1

Putting apps in a sandbox

Apps can only do things that Apple allows and cannot ask (or persuade, or trick) the user for permission to do 'dangerous' things.



- Would this break my phone?
- Would this run my battery down?
- Steal my bank details?

Technical idea #2

Distributing software through a centralized app store

Download Search Apps... Windows Mac iOS Android More

Download offers the opportunity to buy software and apps. When you buy through our links, we may get a commission.

Home / Windows / Utilities & Operating Systems / File Compression / WinRAR (64-bit)

WinRAR (64-bit)

By RARLAB FREE TO TRY

[DOWNLOAD NOW](#) [PREMIUM UPGRADE](#)

Key Details of WinRAR (64-bit)

- Take full control over RAR and ZIP archives, along with unpacking a dozen other archive formats
- Last updated on 11/25/21
- There has been 1 update within the past 6 months
- The current version has [1 flag on VirusTotal](#)
- Also available on [Android](#) and [Mac](#)

Archive name and parameters: Downloads.rar

Downloads (evaluation copy):

Name	Size	Type	Modified
convert		File folder	5/22/2019 10:23 AM
older_downloads		File folder	5/22/2019 9:48 AM
backup_data.rar	10,573,324	WinRAR archive	5/22/2019 10:14 AM
Documents.rar	4,587,398	WinRAR archive	5/22/2019 10:10 AM
Downloads.rar	2,964,465	WinRAR archive	5/22/2019 10:03 AM
music.zip	15,163,290	WinRAR ZIP archive	5/22/2019 10:22 AM

13.51 4G

skype for iphone

Skype for iPhone

Social Networking

★★★★☆ 7

[GET](#)

In-App Purchases

GroupMe

Social Networking

★★★★☆ 56

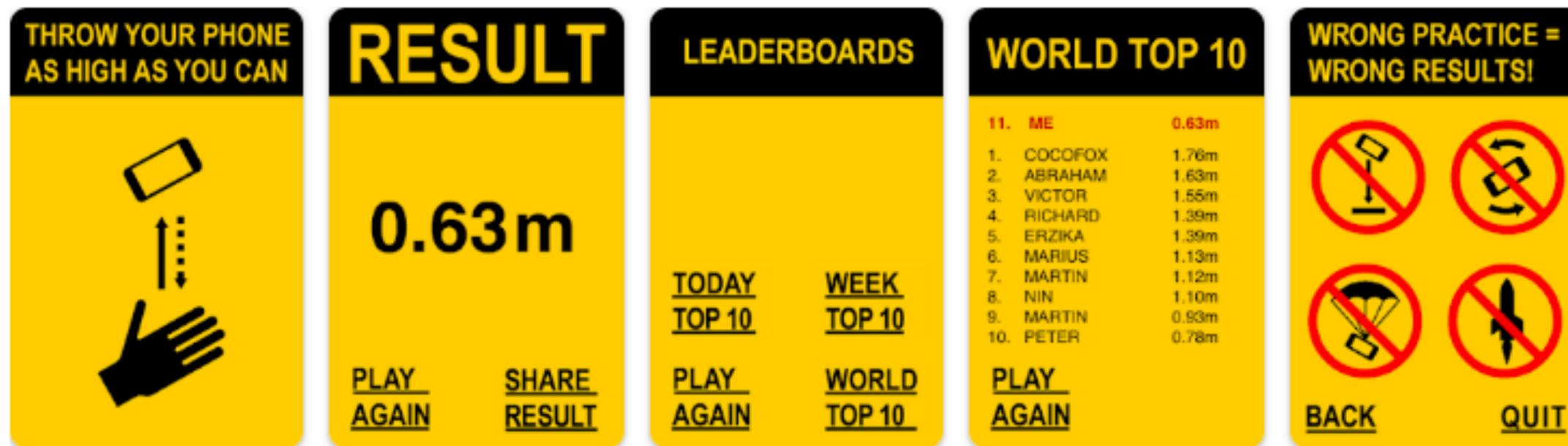
[GET](#)

In-App Purchases

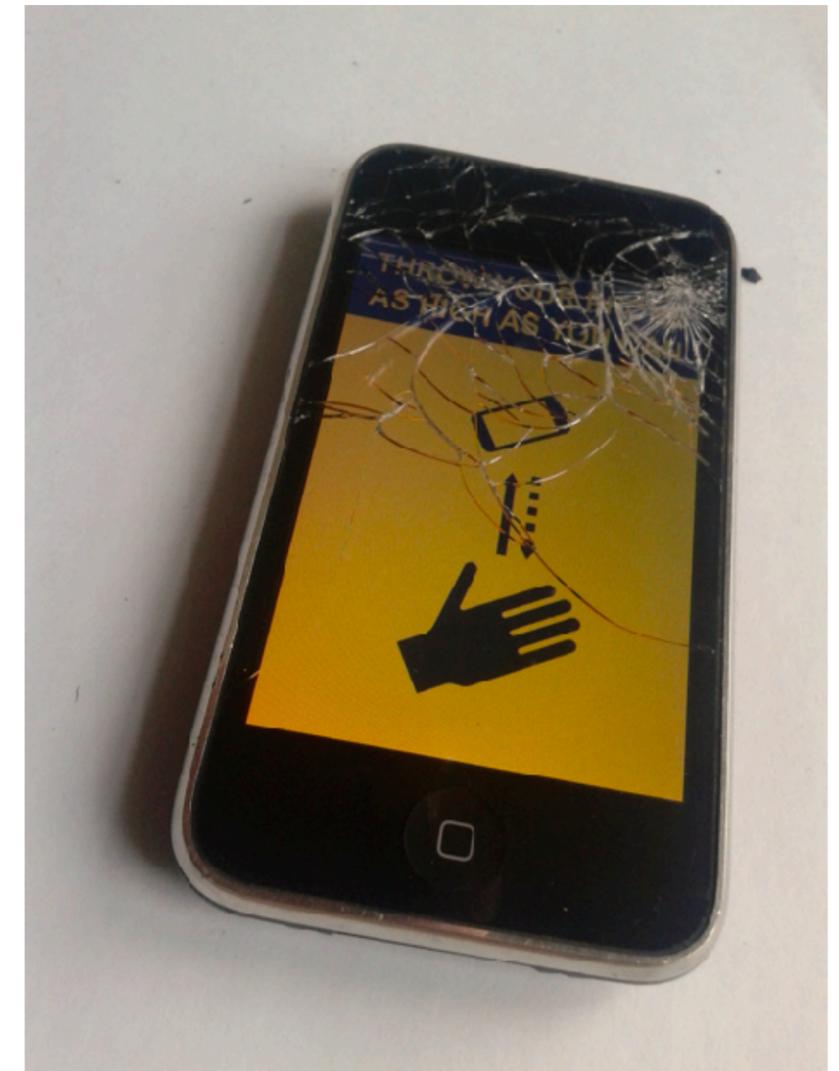
Today Games Apps Updates Search

Technical idea #2

Distributing software through a centralized app store



S.M.T.H.: Send Me to Heaven



Issues around idea #1

Putting apps in a sandbox

Apps can only do things that Apple allows and cannot **ask** (or **persuade**, or **trick**) the user for **permission** to do '**dangerous**' things.

- What are 'dangerous' things?
- How can we trust Apple?
- How can we detect if apps trick the user?
-

Issues around idea #1

After 15 years - still iterating



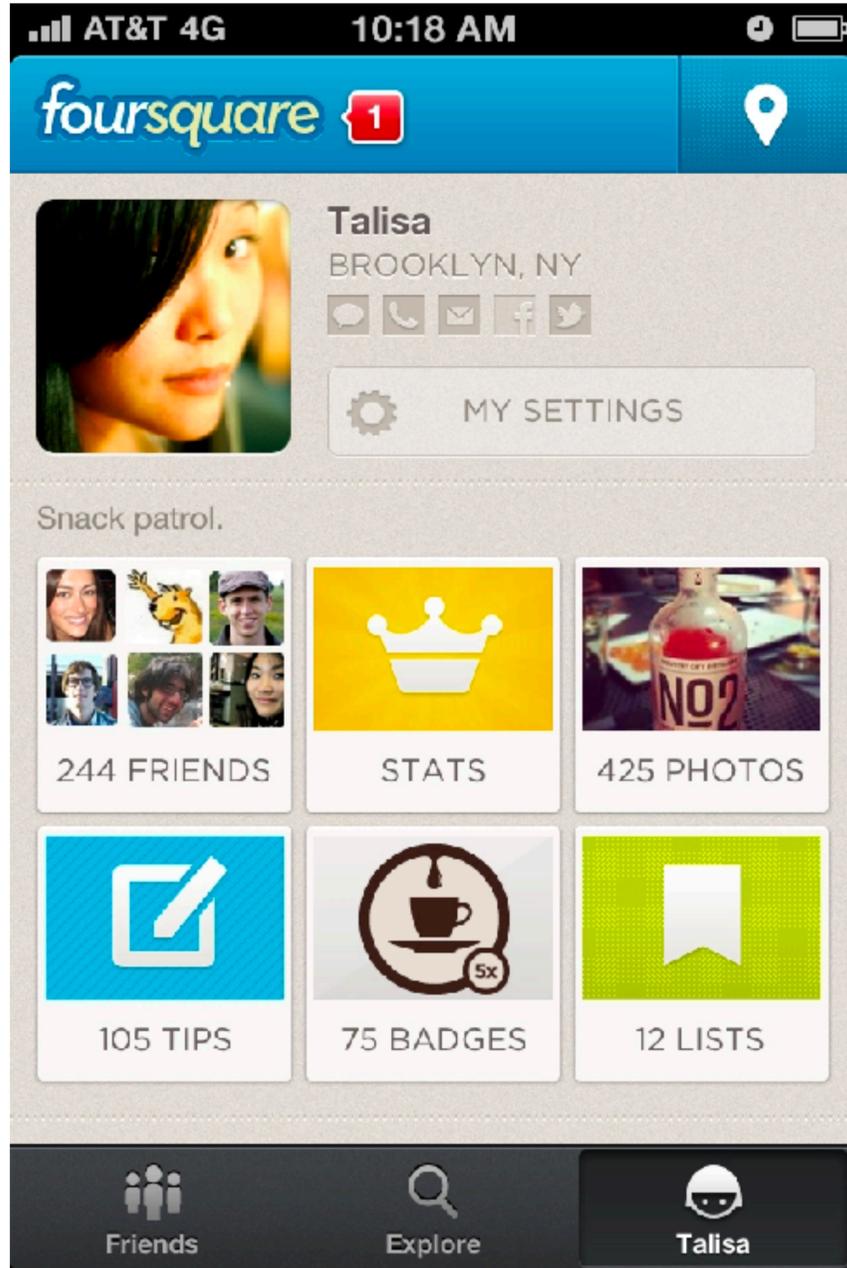
Issues around idea #2

Distributing software through a centralized app store



Issues around idea #2

Software often has a cloud component.

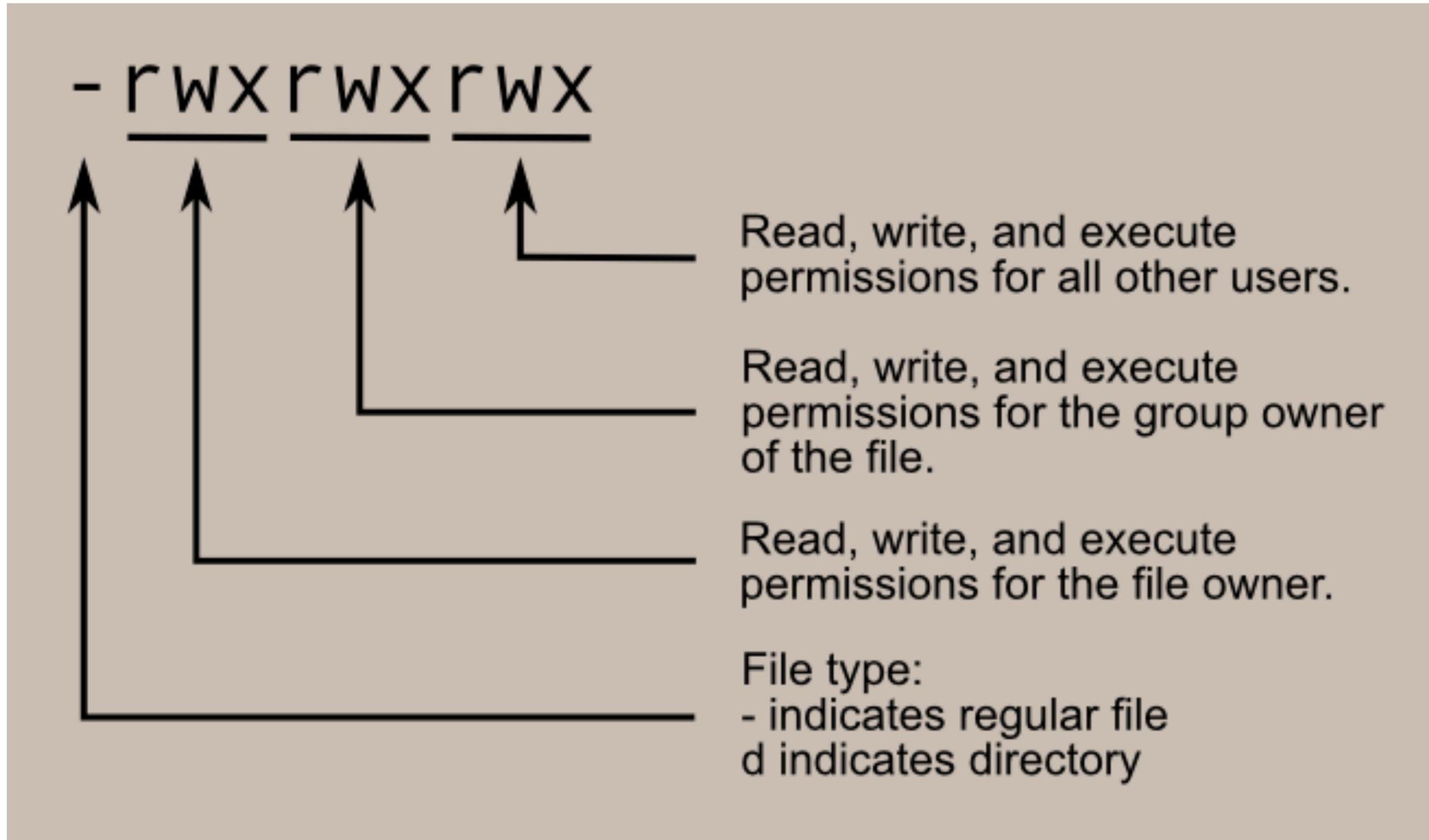


Linux systems

```
aditya314@ubuntu: ~
aditya314@ubuntu:~$ ls
Desktop      examples.desktop  Music      Public      Videos
Documents    ggf.txt           new one    Templates    xyz.txt
Downloads    listfile          Pictures    Untitled Document
aditya314@ubuntu:~$ ls -l
total 52
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Desktop
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Documents
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Downloads
-rw-r--r-- 1 aditya314 aditya314 8980 Mar  5 01:05 examples.desktop
-rw-rw-r-- 1 aditya314 aditya314   0 Mar  5 03:53 ggf.txt
-rw-rw-r-- 1 aditya314 aditya314   0 Apr 27 02:47 listfile
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Music
drwxrwxr-x 2 aditya314 aditya314 4096 Mar  5 03:53 new one
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Pictures
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Public
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Templates
-rw-rw-r-- 1 aditya314 aditya314   0 Apr 27 02:55 Untitled Document
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Videos
-rw-rw-r-- 1 aditya314 aditya314  268 Mar  5 04:17 xyz.txt
aditya314@ubuntu:~$
```

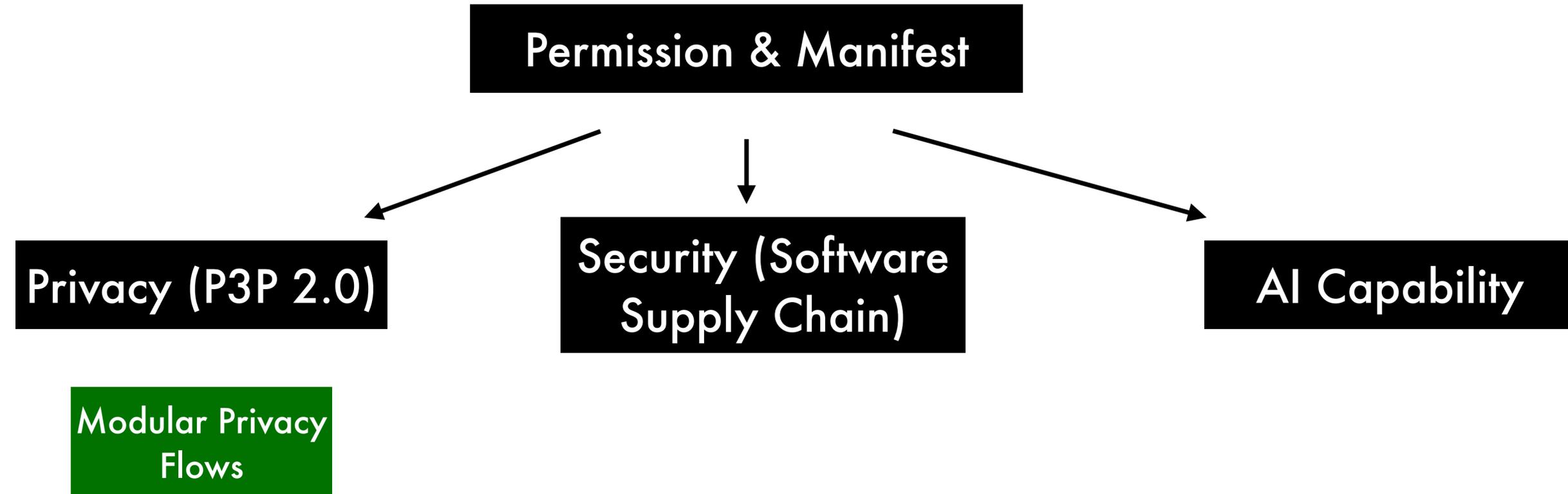
- Users
- File permissions

File permissions

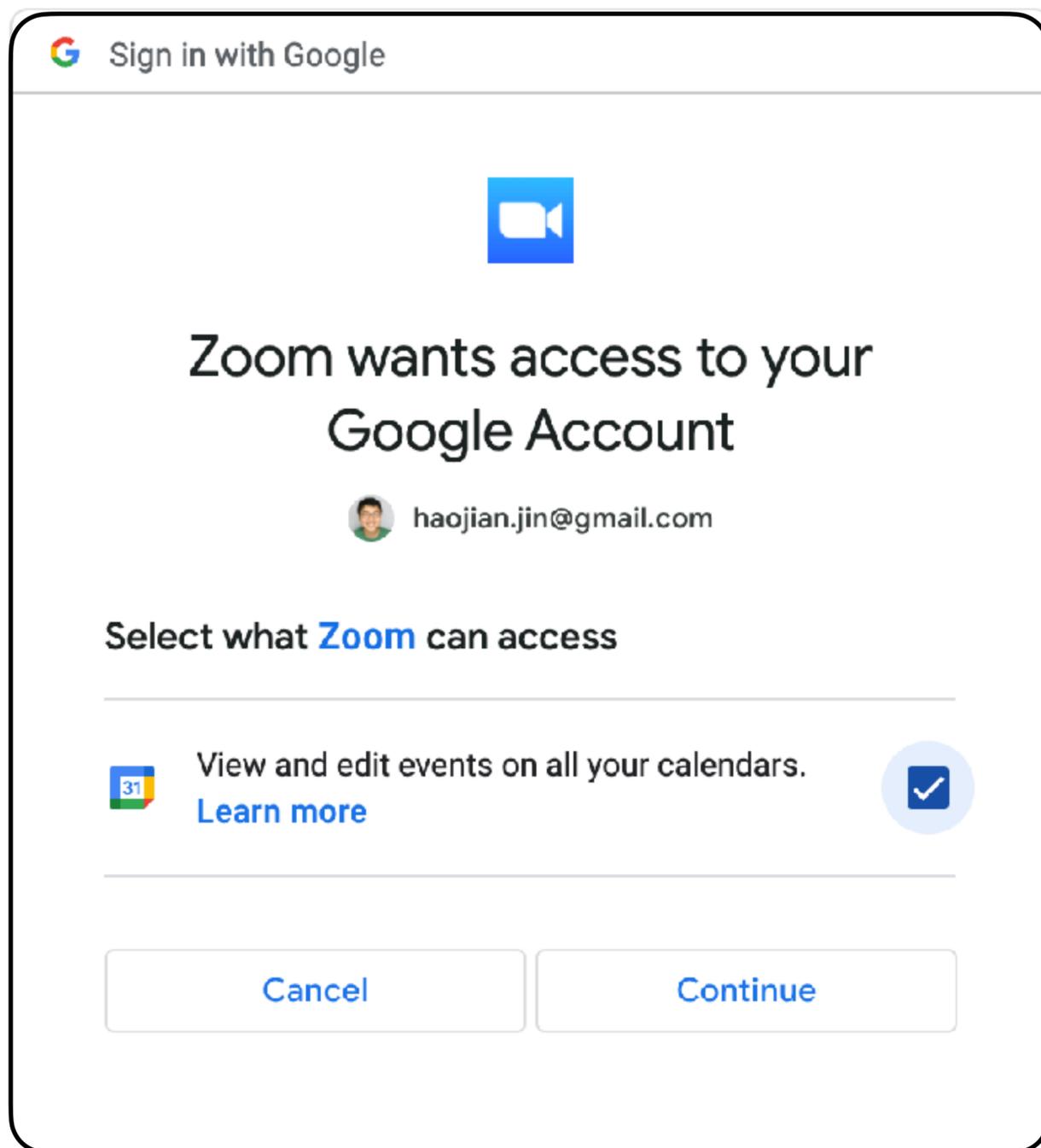


Current research focus

Next generation permission system

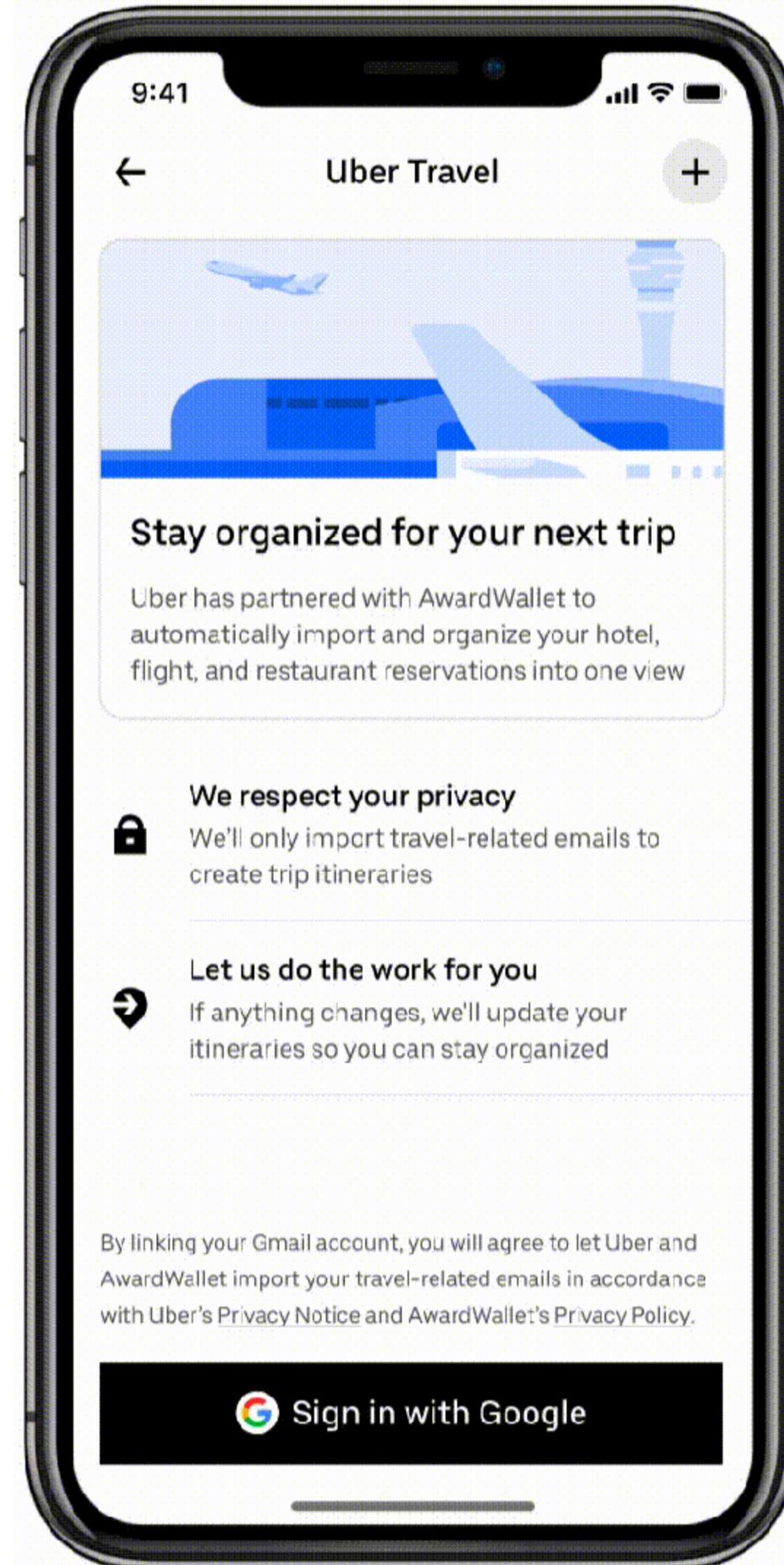


Zoom accesses **all** your calendar events **continuously**!



Calendar events that contain
<https://zoom.us/xxxxxx>

Uber wants to see
all your emails.



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)

Principle of least privilege

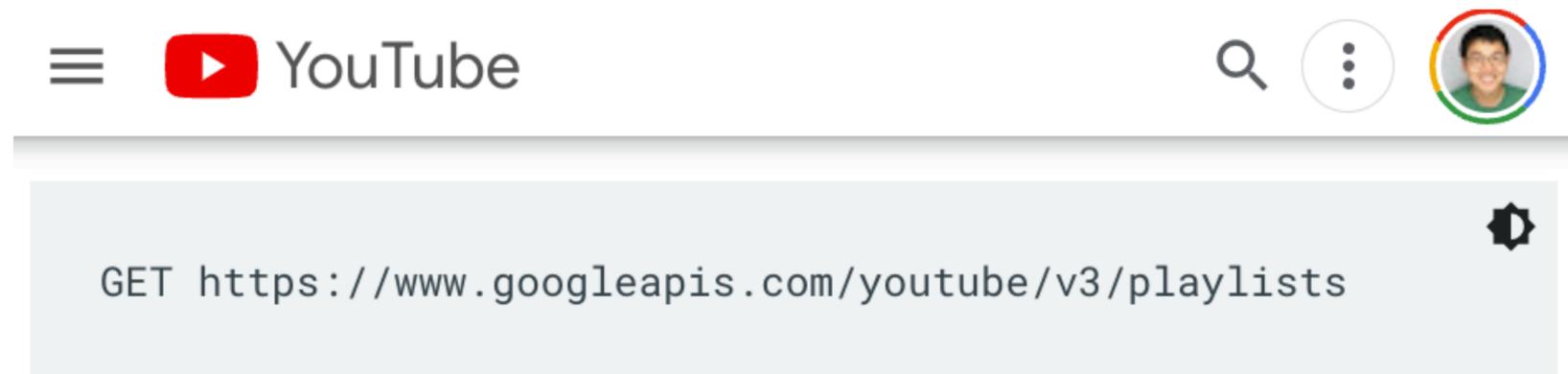
*"A security architecture should be designed so that each entity is granted the **minimum system resources** and authorizations that the entity needs to perform its function.."*

Google APIs - All-or-nothing binary permissions

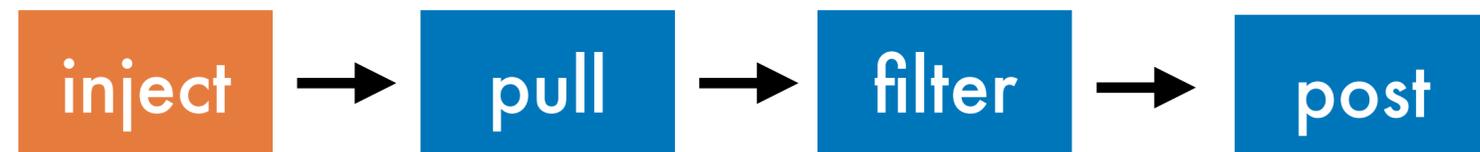
Scope	Meaning
<code>https://www.googleapis.com/auth/calendar</code>	read/write access to Calendars
<code>https://www.googleapis.com/auth/calendar.readonly</code>	read-only access to Calendars
<code>https://www.googleapis.com/auth/calendar.events</code>	read/write access to Events
<code>https://www.googleapis.com/auth/calendar.events.readonly</code>	read-only access to Events
<code>https://www.googleapis.com/auth/calendar.settings.readonly</code>	read-only access to Settings
<code>https://www.googleapis.com/auth/calendar.addons.execute</code>	run as a Calendar add-on

Program data transformation functions using chainable *operators*

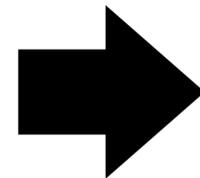
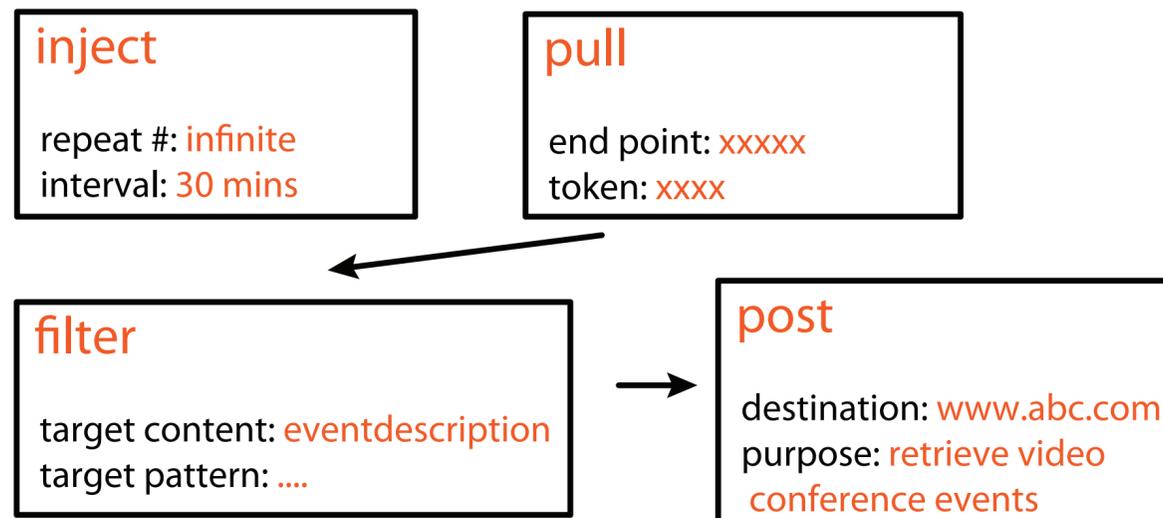
URL-based APIs



Operator-based APIs

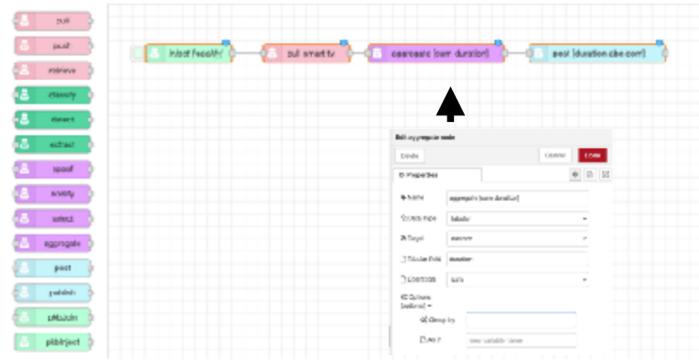


A text-based whitelist *manifest* (i.e., program representation)



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

System builders



Offer a set of operators as
the API



Execute the manifest using
preloaded implementations

Developers

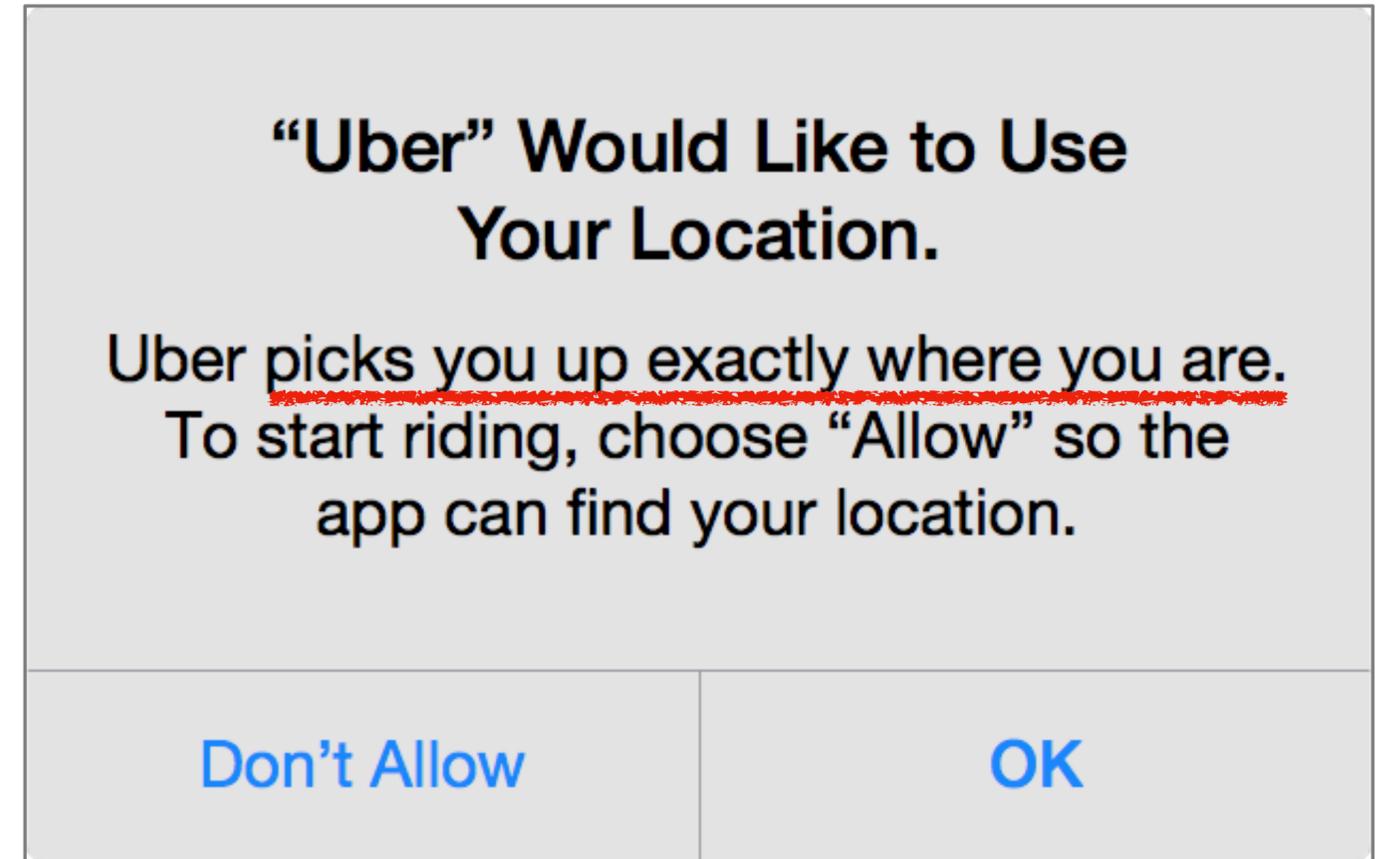
```
@purpose: The app can access calendar events  
which contains a zoom link.  
ZoomCalendarIntegration{  
  // operator [properties]  
  inject[...] -> pull Calendar[...] ->  
  filter [Zoom join link pattern] ->  
  post [Zoom events]  
}
```

Author a manifest by
connecting operators

Talk outline

1. Modular Privacy Flows (MPF) in a Nutshell
2. **Why** MPF
3. **How** MPF
4. **When and when not** MPF
5. Future Work

Purpose strings.

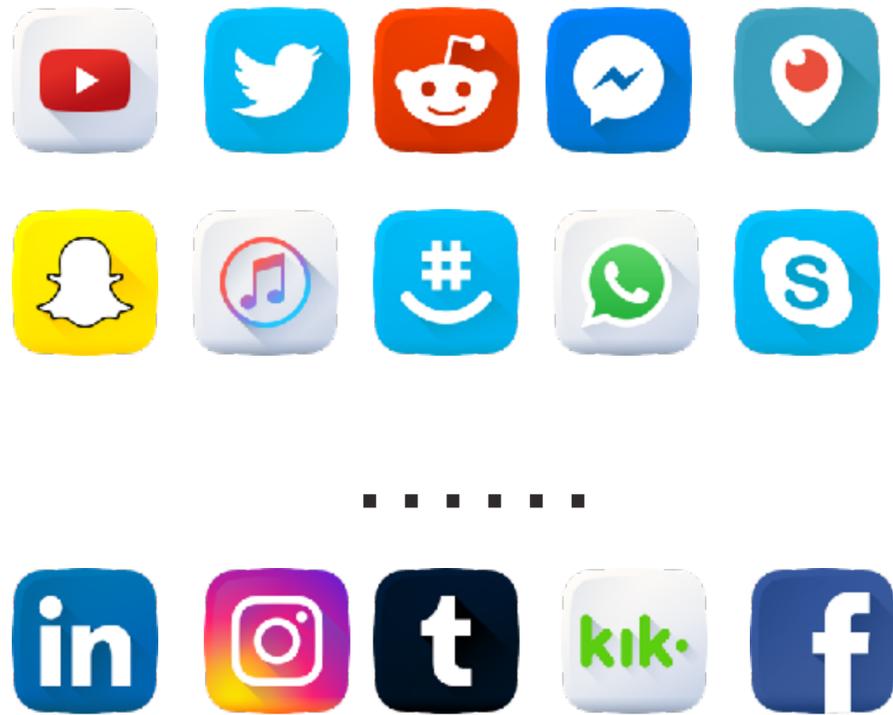


arbitrary text,
manually annotated,
hard to validate/assess.



MobiPurpose is a scalable in-lab solution that can index fine-grained privacy attributes (who, where, what, why) of outgoing network requests.

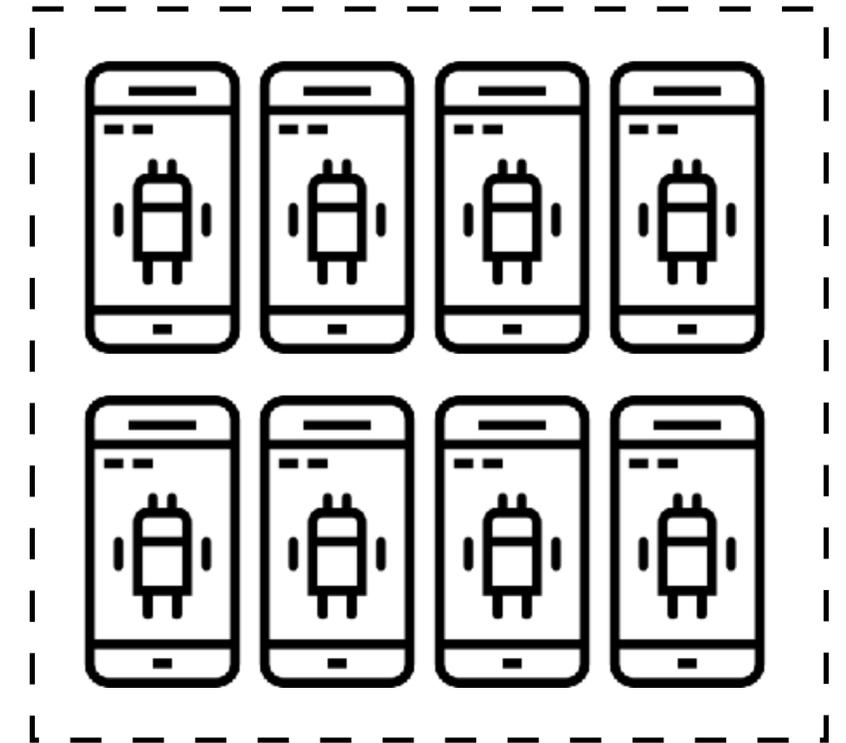
MobiPurpose - network tracing



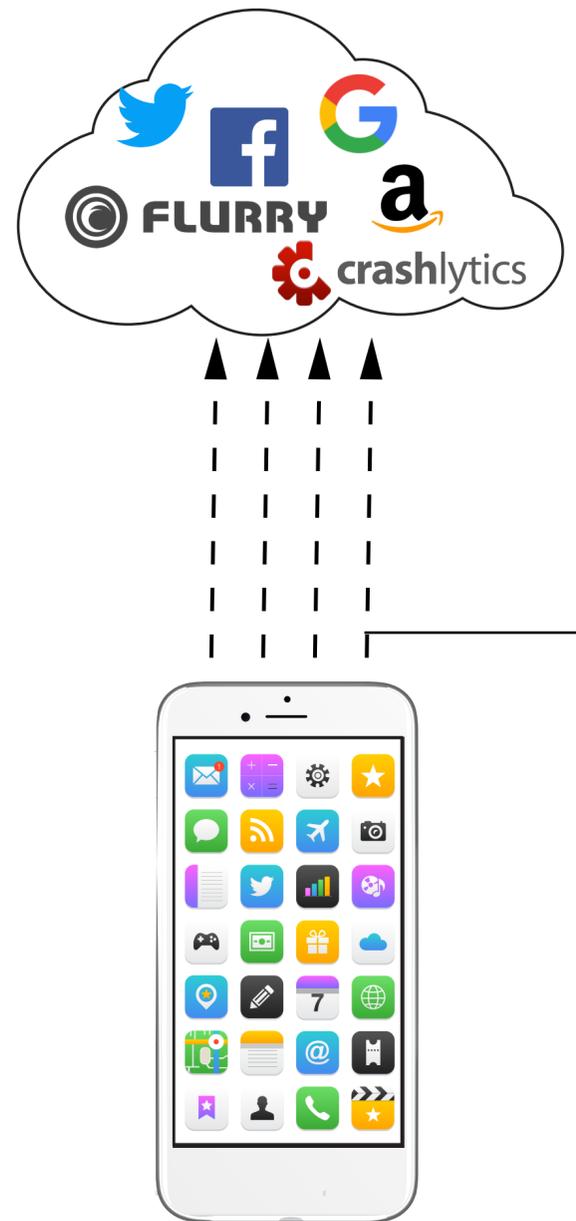
185k apps



Intercepted 2 million
unique traffic requests



MobiPurpose - network tracing



Traffic request snapshot

source app:

com.inkcreature.predatorfree

Who?

connect to host:

inkcreature.com

Where?

server path:

/_predatorServer/

key-value pairs in request body:

myLat: 40.4435877

myLon: -79.9452883

....

Key-value pairs

Traffic request snapshot

```
source app:  
  com.inkcreature.predatorfree  
connect to host:  
  inkcreature.com  
server path:  
  /_predatorServer/  
  
key-value pairs in request body:  
  myLat: 40.4435877  
  myLon: -79.9452883  
  ....
```

2,008,912 unique traffic requests
from 14,910 apps

contacting

12,046 unique domains

302,893 unique URLs

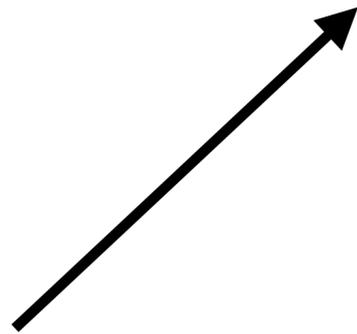
We publish the dataset at:
<http://bit.ly/purposedata>

DATA TYPES

DATA PURPOSES

EXAMPLES

location



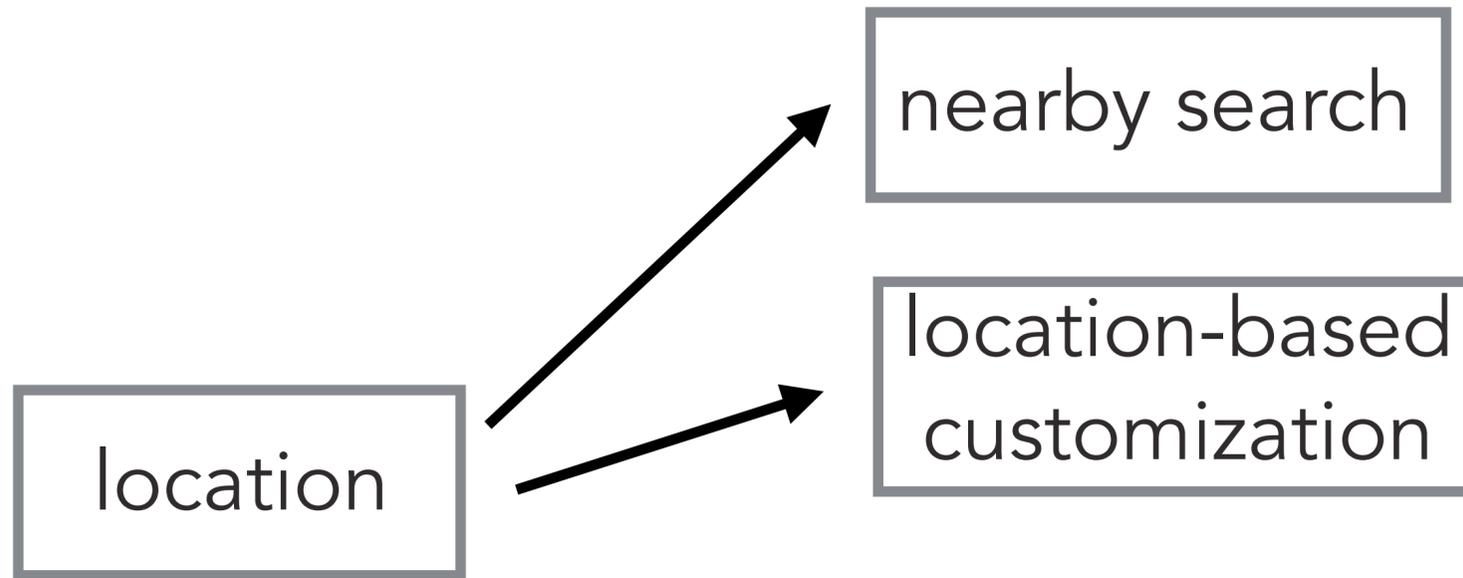
nearby search



DATA TYPES

DATA PURPOSES

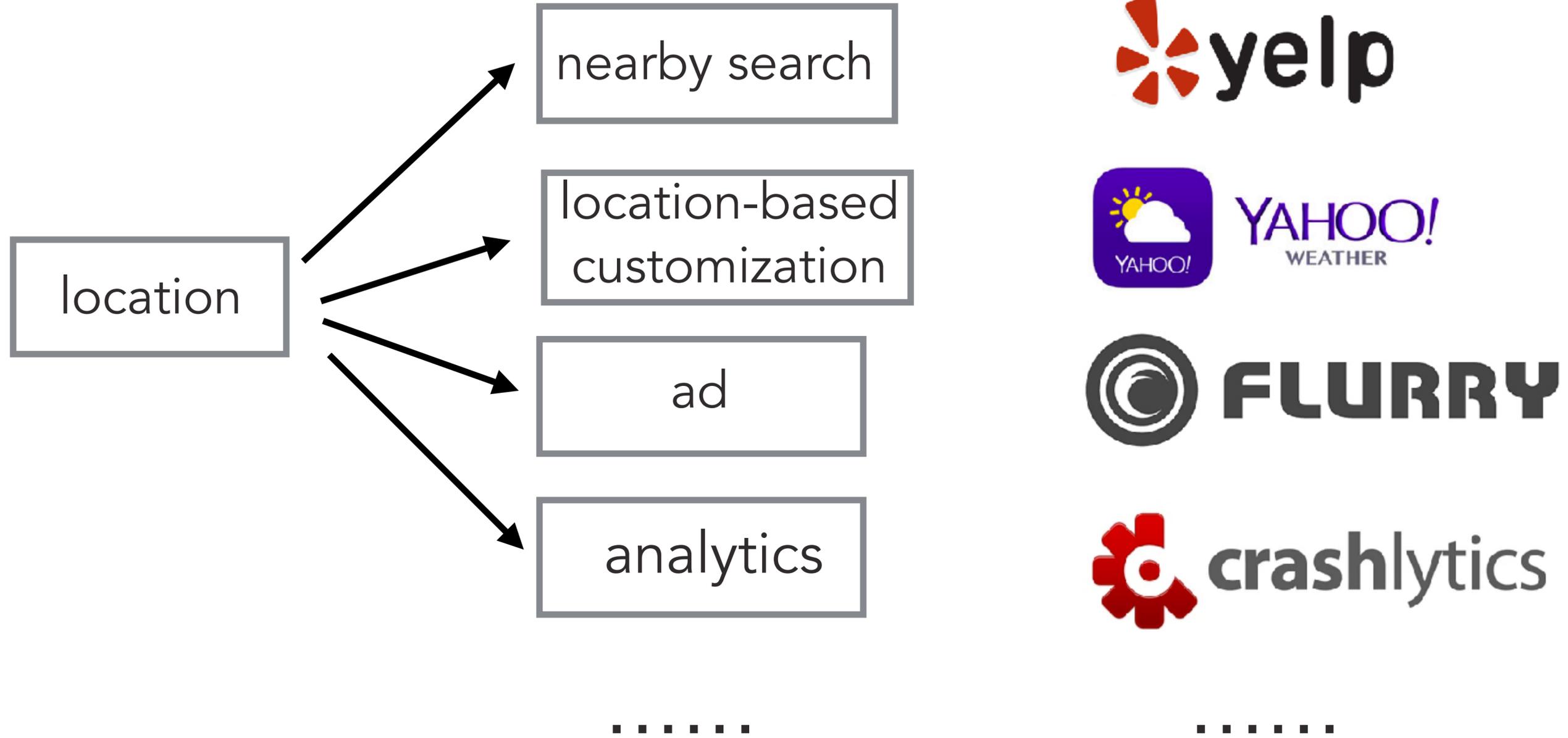
EXAMPLES



DATA TYPES

DATA PURPOSES

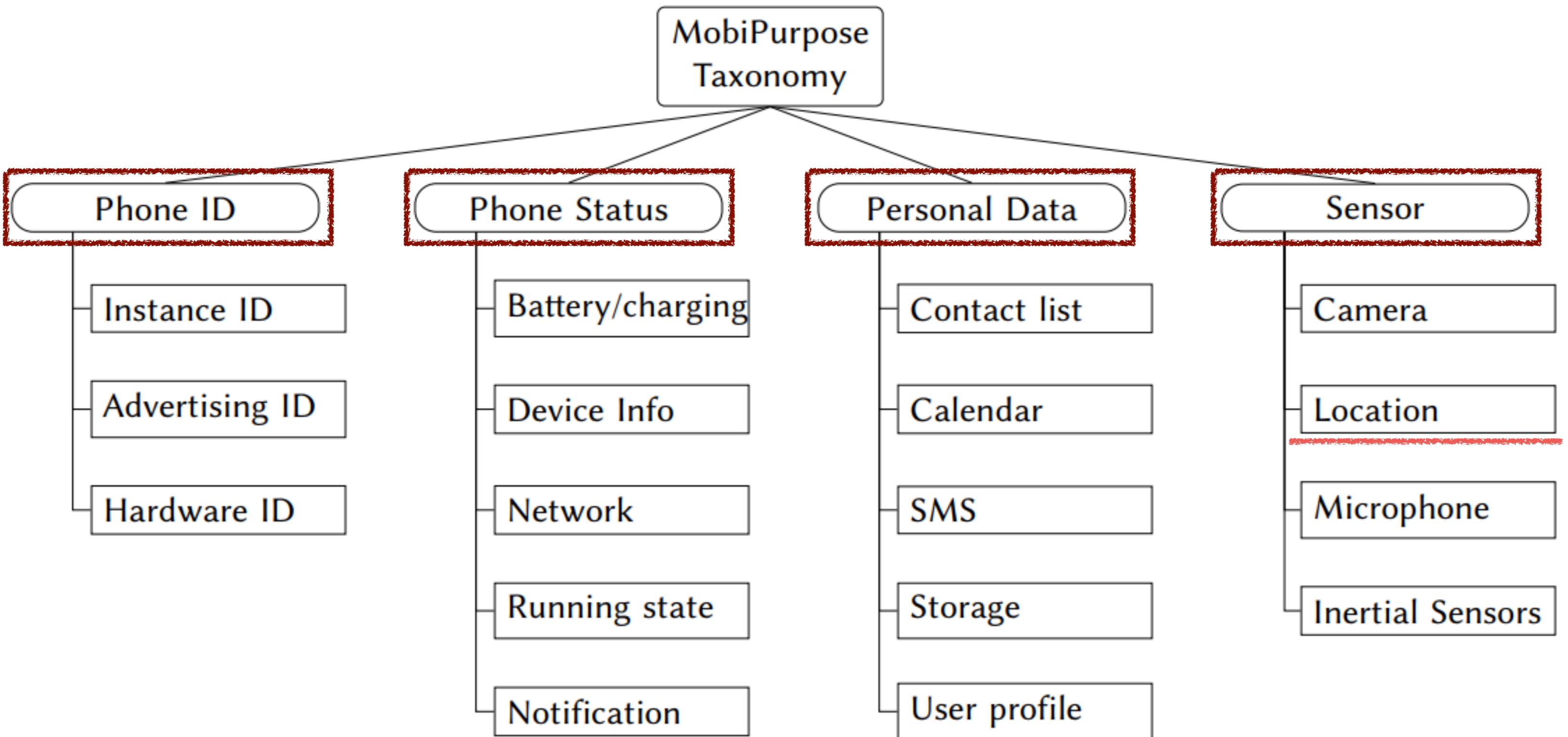
EXAMPLES



13 common data collection purposes for location data

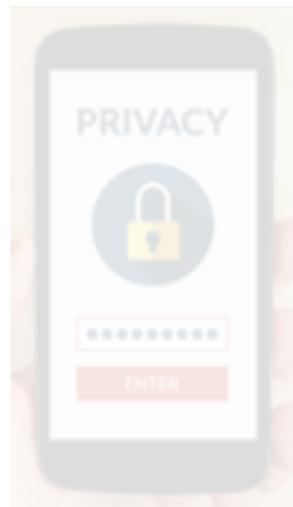
Location ⁷	Nearby Search	Search nearby POIs/real estates
	Location-based Customization	Fetch local weather/radio information
	Query Transportation Information	Estimate the trip time through Uber API
	Recording	Track the running velocity
	Map and Navigation	Find the user location in Map apps
	Geosocial Networking	Find nearby users in the social network
	Geotagging	Tag photos with locations
	Location Spoofing	Set up fake GPS locations
	Alert and Remind	Remind location-based tasks
	Location-based game	Play games require users' physical location
	Reverse geocoding	Use the GPS coords to find the real world address.
	Data collection for analytics	Collect data for marketing analysis
	Data collection for ad	Collect data for ad personalization

See the complete taxonomy at:
<http://bit.ly/mobitaxonomy>



data types

Three empirical studies



Large-scale mobile
network tracing^[1]



Smart home^[2]
applications



Smart city^[3]
applications

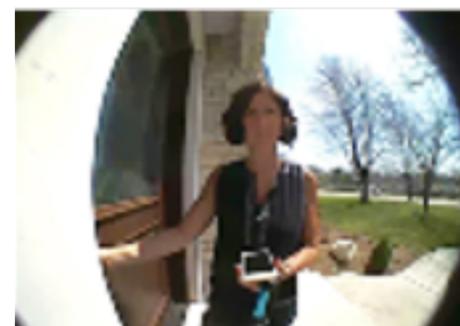
77% Smart home apps do not need raw data.

Sensor

Raw

Needed data

Hello visitor

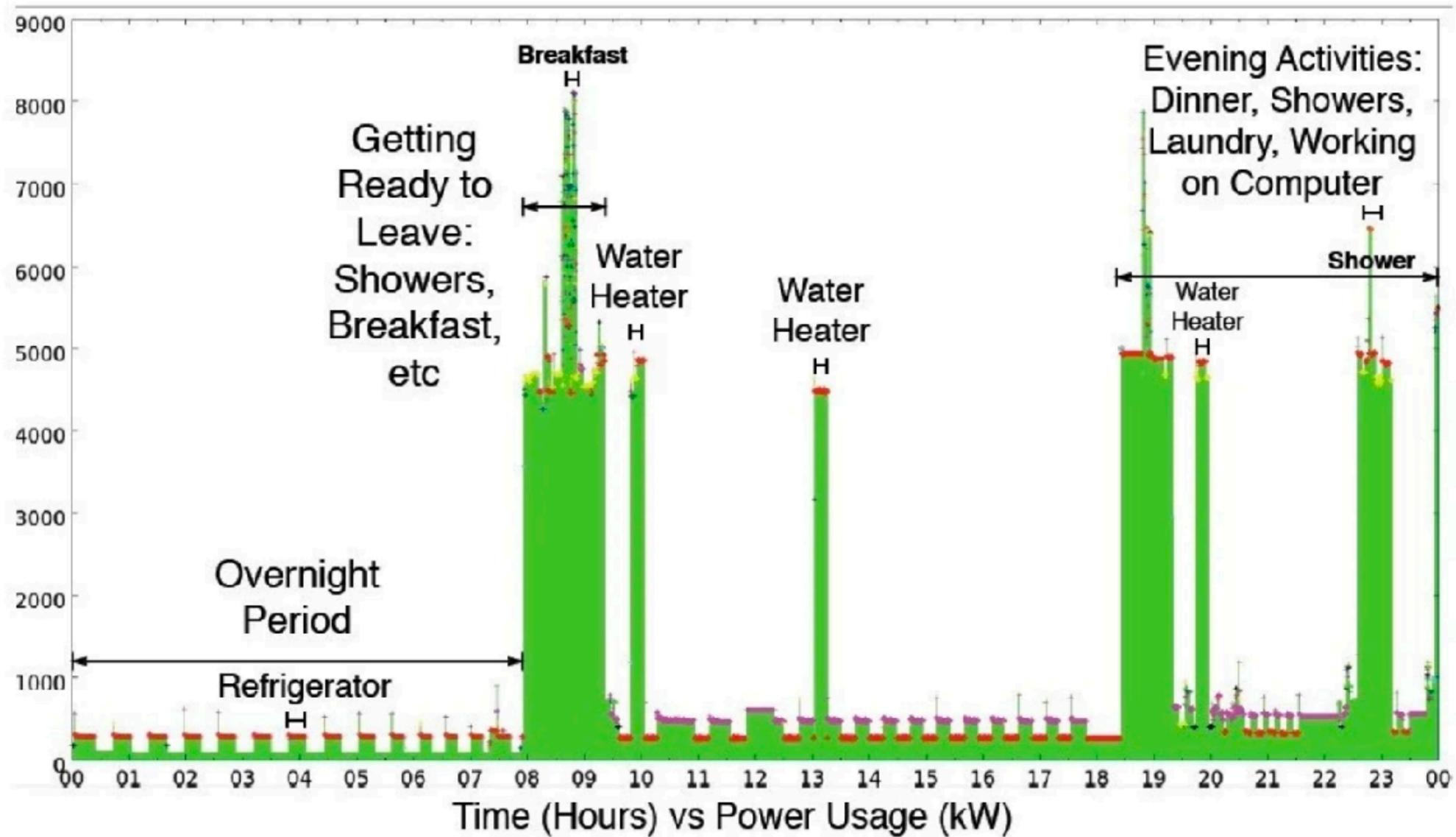


Noise level

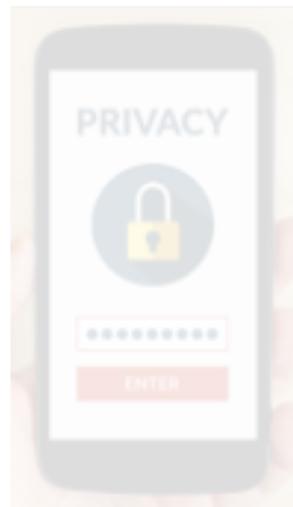


55 db

Oversensing



Three empirical studies



Large-scale mobile
network tracing^[1]



Smart home^[2]
applications



Smart city^[3]
applications

76/80 Smart city apps only need aggregated data.

Aggregation types	Applications
Count/Sum/Average	Electricity usage (sum/average), water consumption (sum/average), heat consumption (sum/average), water leaks (sum/average), fire alerts (average), CO2 emission tracking (sum), power-efficient appliances (sum/percentage), indoor vs. outdoor water usage (group sum), building vs. home energy efficiency (group sum), electric vehicle usage (sum), waste management (sum), renewable energy (sum), rainwater collection (sum), public transport tracking (count), intersection traffic (count), residents physical activity (average), pedestrian entrance tracking (count), pedestrian foot traffic (count), user mobility (count), time spent asleep (average), number of users that pass by a billboard (count), DMV visitor counting (GPS), fire department delay (average), traffic light wait time (average), irrigation tracking (sum/average), electromagnetic field level monitoring (sum/average), food consumption (sum/average), indoor home temperature (average), parking occupancy (count/average/sum)
Rank/Median/Top X	Snow plowing fairness, public space noise level monitoring
Histogram/Heatmap	Network speed monitoring (group average), water quality monitoring (group average), food consumption (group average), street light brightness (group average), noisiest neighborhoods (group average)
Route	Delivery route with best parking availability, Smart garbage pickup route, least-polluted biking path
TF-IDF	Determine restaurant type for region, trending TV search queries
Trilateration	Gunshot localization, air pollution source localization
Time-series forecasting	Parking occupancy prediction
Non-aggregatable apps	Searching for full trash cart, On-demand rides with autonomous vehicles, DMV visitors Counting through a camera, Garbage collection prediction

New understanding about privacy

State of the art

25% users were surprised this app sent their **approximate location** to dictionary.com for searching nearby words.

arbitrary text
manually annotated
hard to validate/assess



Sensor

Raw

Needed data

Hello visitor



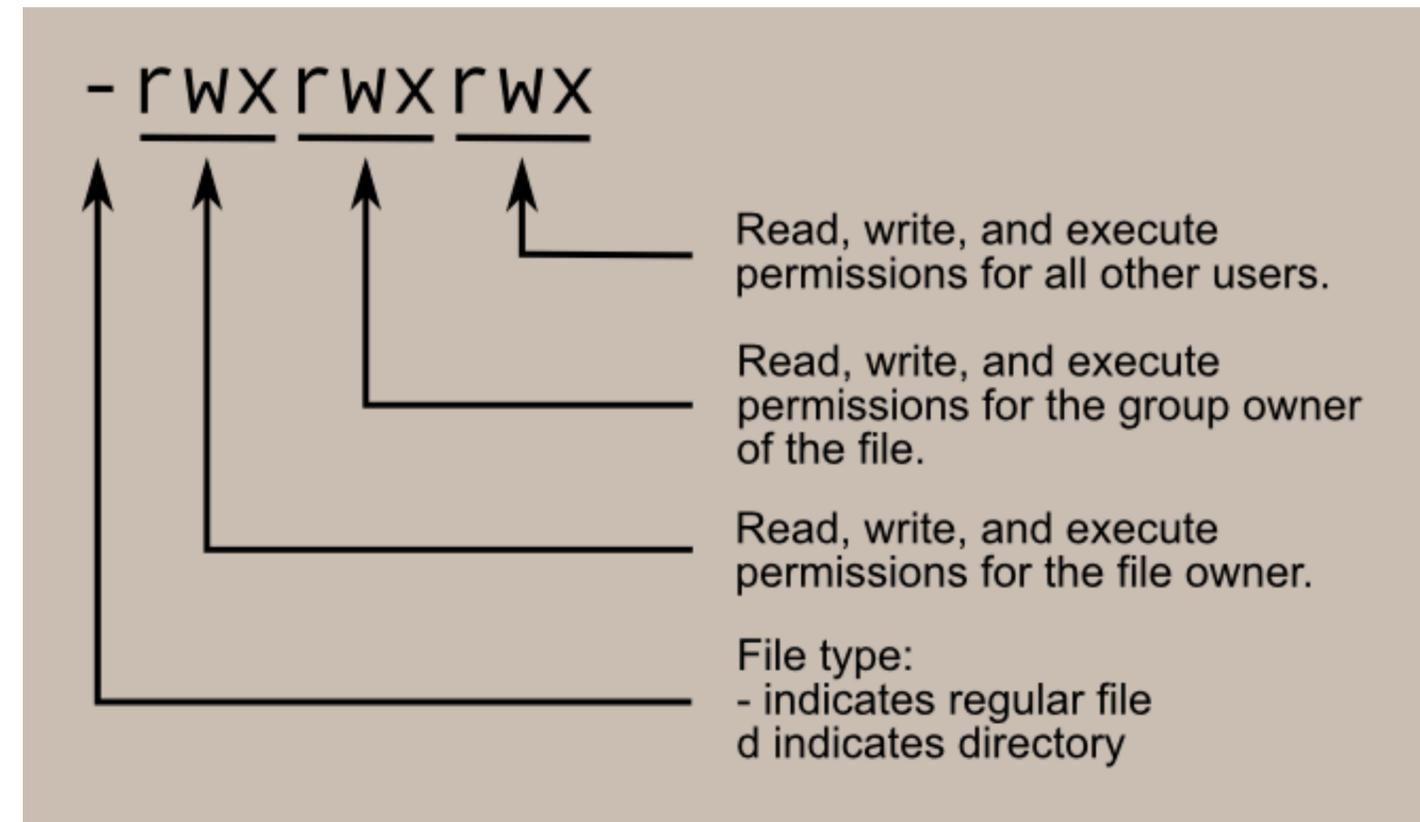
MPF

Enumerable data collection purposes.

Given a purpose, developers do not need raw data.

File permissions are insufficient.

```
aditya314@ubuntu: ~  
aditya314@ubuntu:~$ ls  
Desktop      examples.desktop  Music      Public      Videos  
Documents    ggf.txt           new one    Templates    xyz.txt  
Downloads    listfile          Pictures    Untitled Document  
aditya314@ubuntu:~$ ls -l  
total 52  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Desktop  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Documents  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Downloads  
-rw-r--r-- 1 aditya314 aditya314 8980 Mar  5 01:05 examples.desktop  
-rw-rw-r-- 1 aditya314 aditya314   0 Mar  5 03:53 ggf.txt  
-rw-rw-r-- 1 aditya314 aditya314   0 Apr 27 02:47 listfile  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Music  
drwxrwxr-x 2 aditya314 aditya314 4096 Mar  5 03:53 new one  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Pictures  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Public  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Templates  
-rw-rw-r-- 1 aditya314 aditya314   0 Apr 27 02:55 Untitled Document  
drwxr-xr-x 2 aditya314 aditya314 4096 Mar  5 01:21 Videos  
-rw-rw-r-- 1 aditya314 aditya314 268 Mar  5 04:17 xyz.txt  
aditya314@ubuntu:~$
```

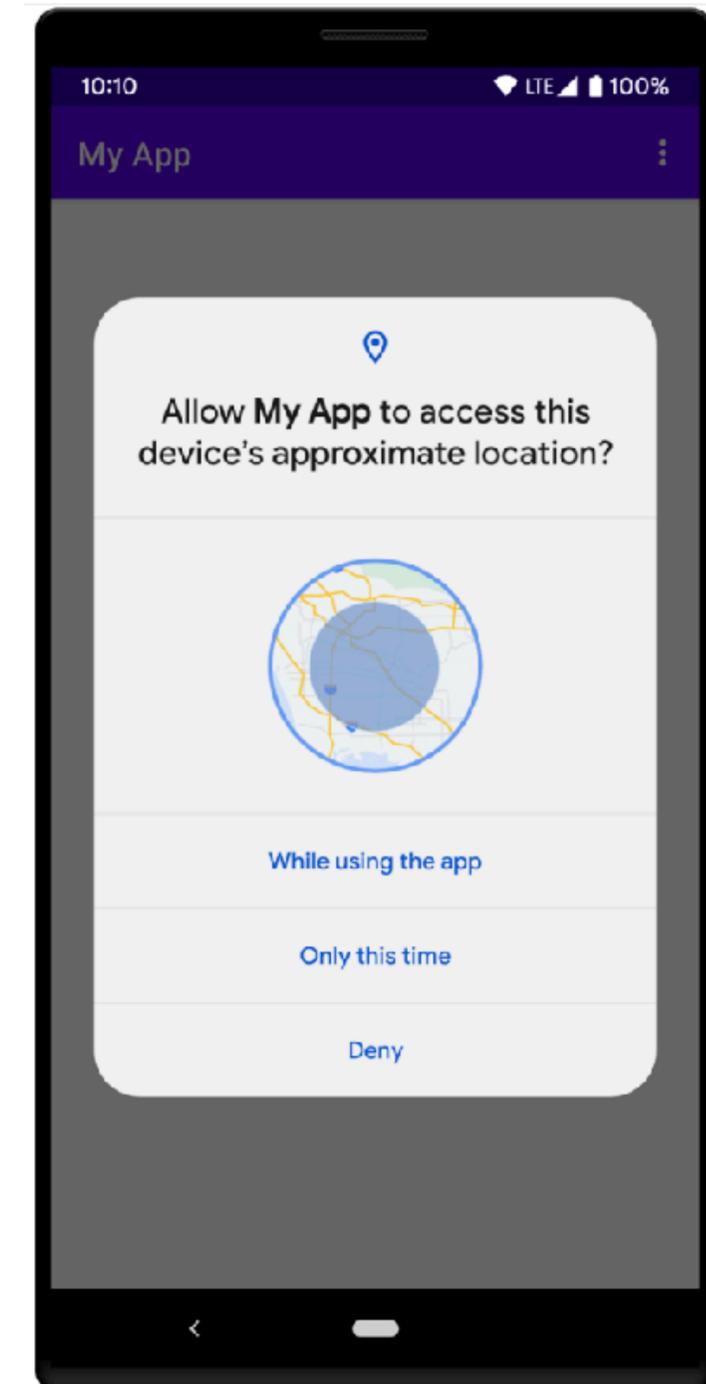
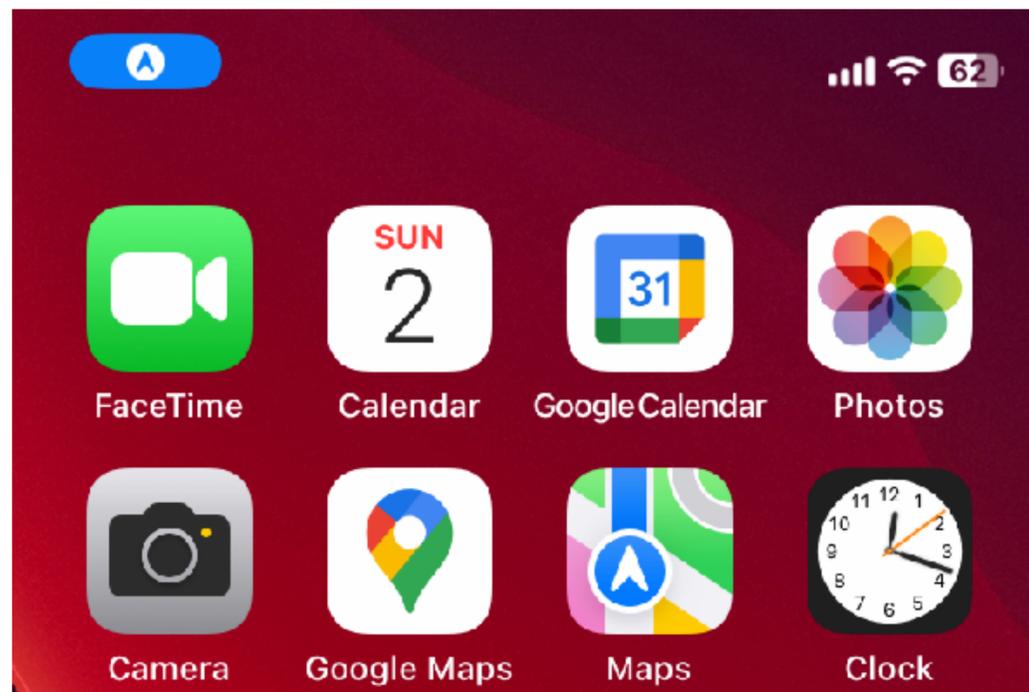


State of the Art: **fine-grained** permission manifest

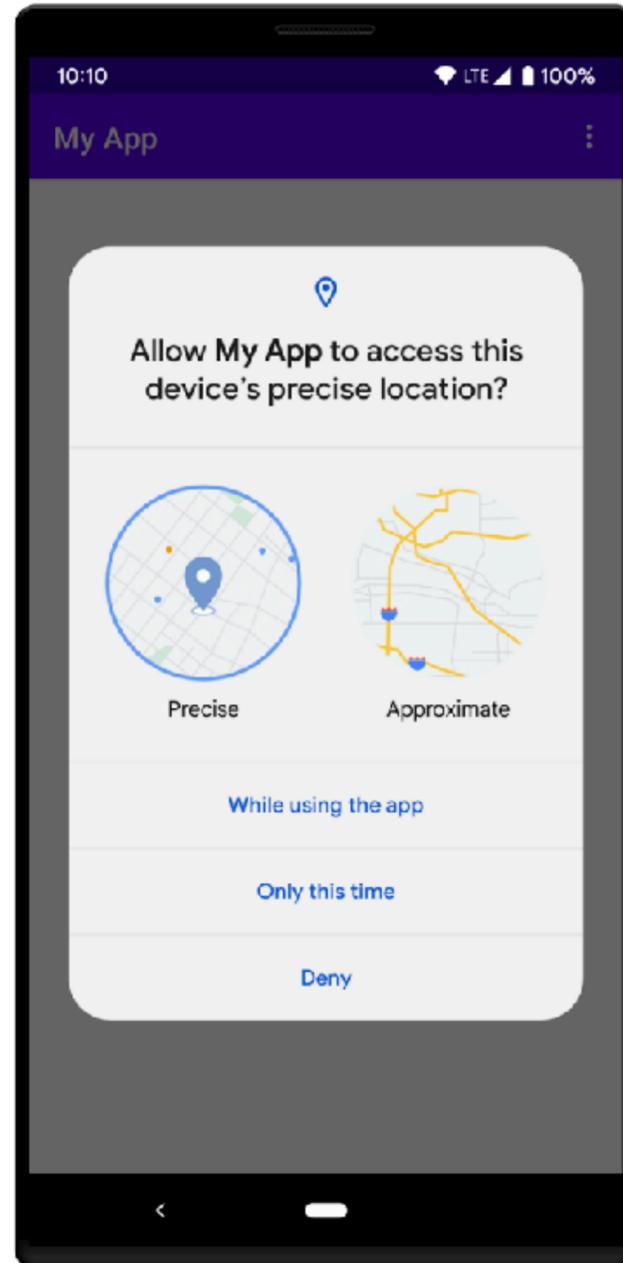
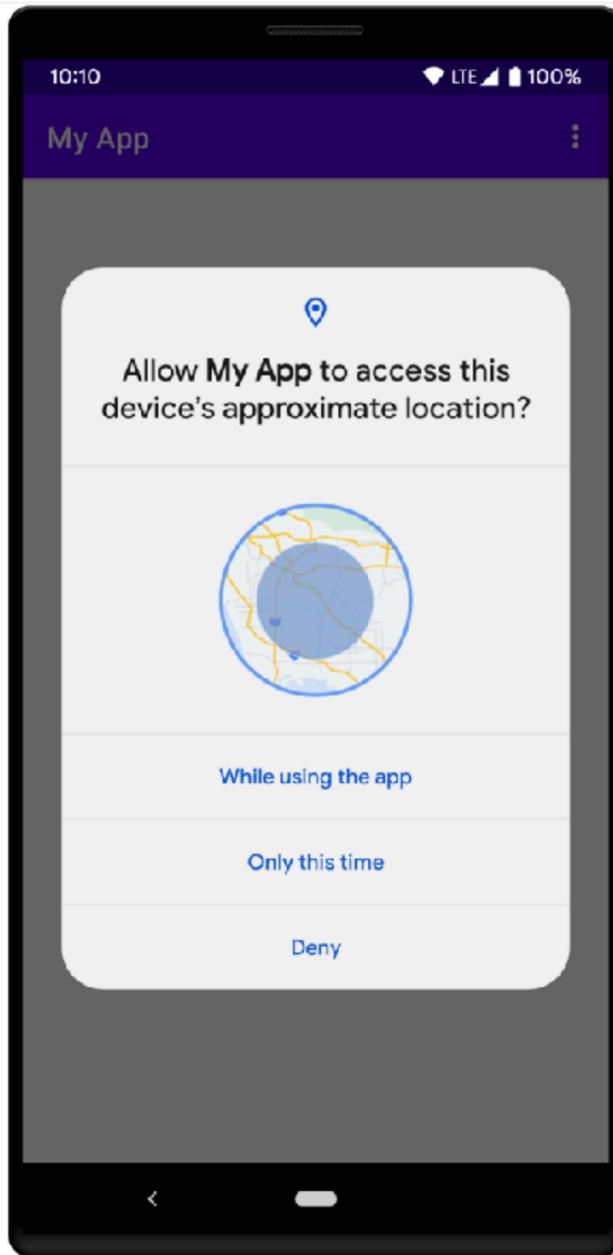
```
<manifest ... >
  <!-- Always include this permission -->
  <uses-permission android:name="android.permission.ACCESS_COARSE_LOCATION" />

  <!-- Include only if your app benefits from precise location access. -->
  <uses-permission android:name="android.permission.ACCESS_FINE_LOCATION" />
</manifest>
```

```
<manifest ... >
  <!-- Required only when requesting background location access on
       Android 10 (API level 29) and higher. -->
  <uses-permission android:name="android.permission.ACCESS_BACKGROUND_LOCATION" />
</manifest>
```



State of the Art: **User choices**



	Precise	Approximate
While using the app	ACCESS_FINE_LOCATION and ACCESS_COARSE_LOCATION	ACCESS_COARSE_LOCATION
Only this time	ACCESS_FINE_LOCATION and ACCESS_COARSE_LOCATION	ACCESS_COARSE_LOCATION
Deny	No location permissions	No location permissions

The permission granularity dilemma

More fine-grained permissions

→ Better privacy

→ More management burden for users

Harder learning curve for app developers

More implementation efforts for system builders

More coarse-grained permissions

→ Worse privacy

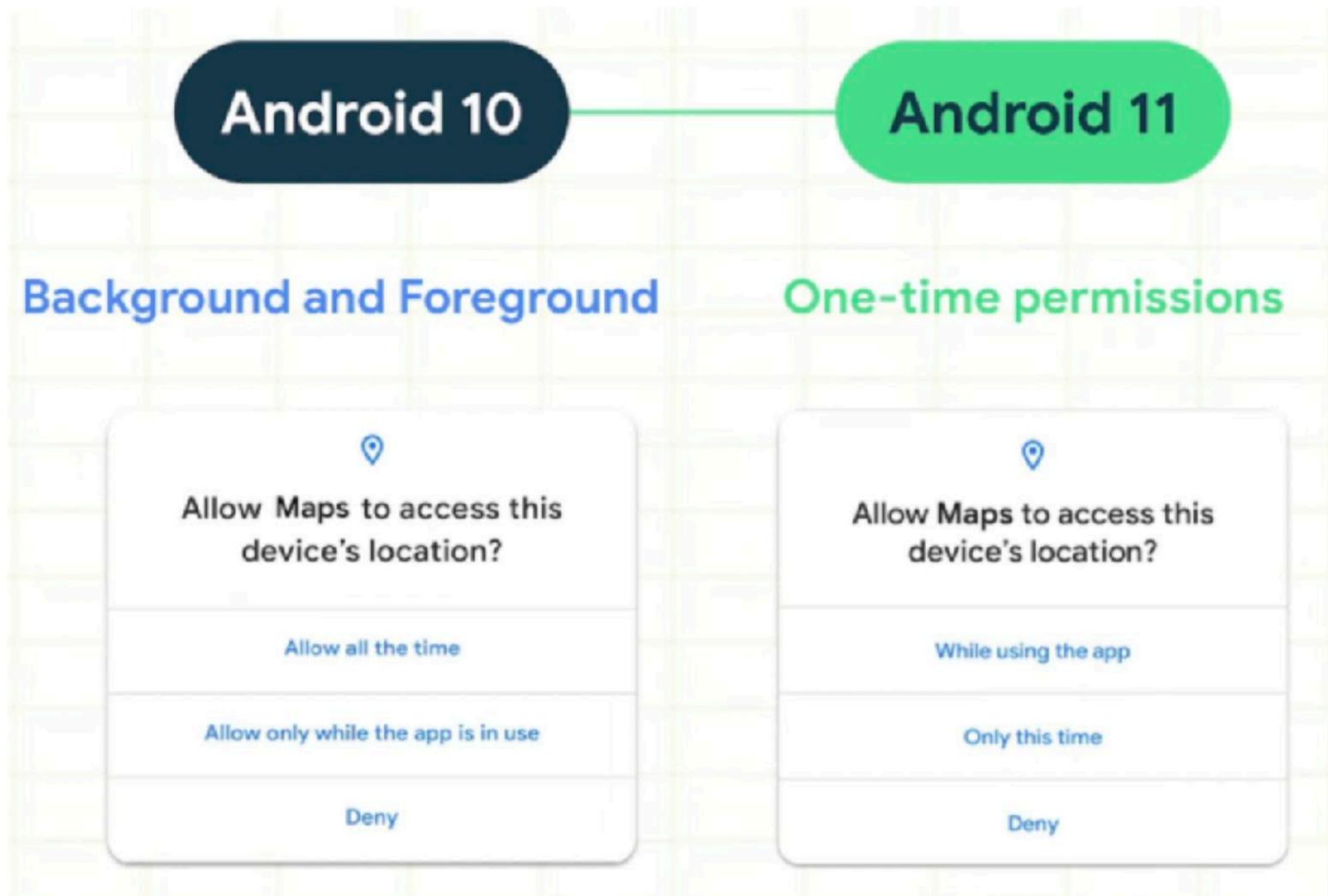
→ Oversensing risks

More users deny data requests

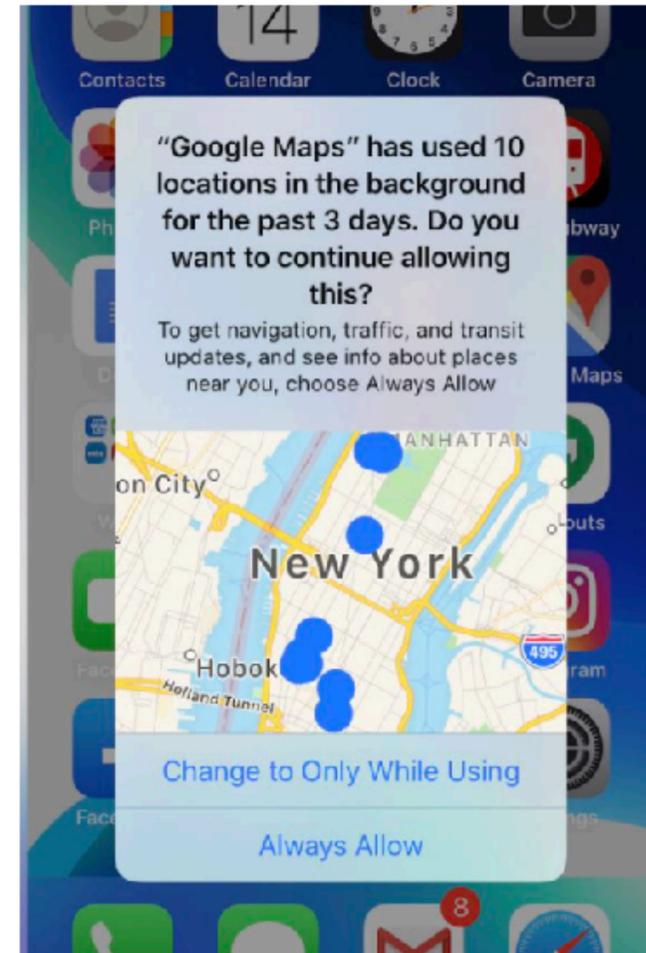
More complaints for system builders

Hard to gain trust from users for app developers

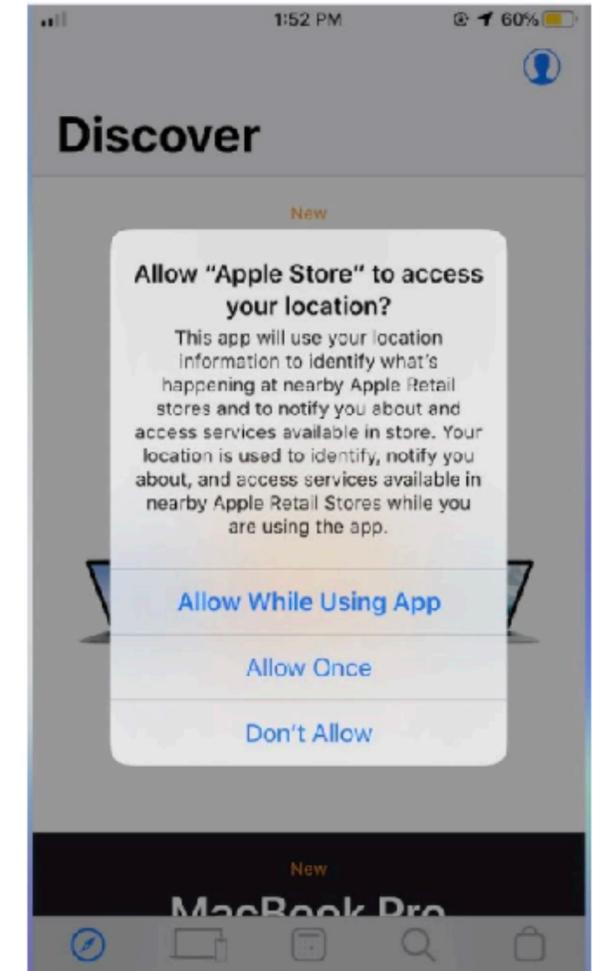
On-going permission dilemma



iOS 12



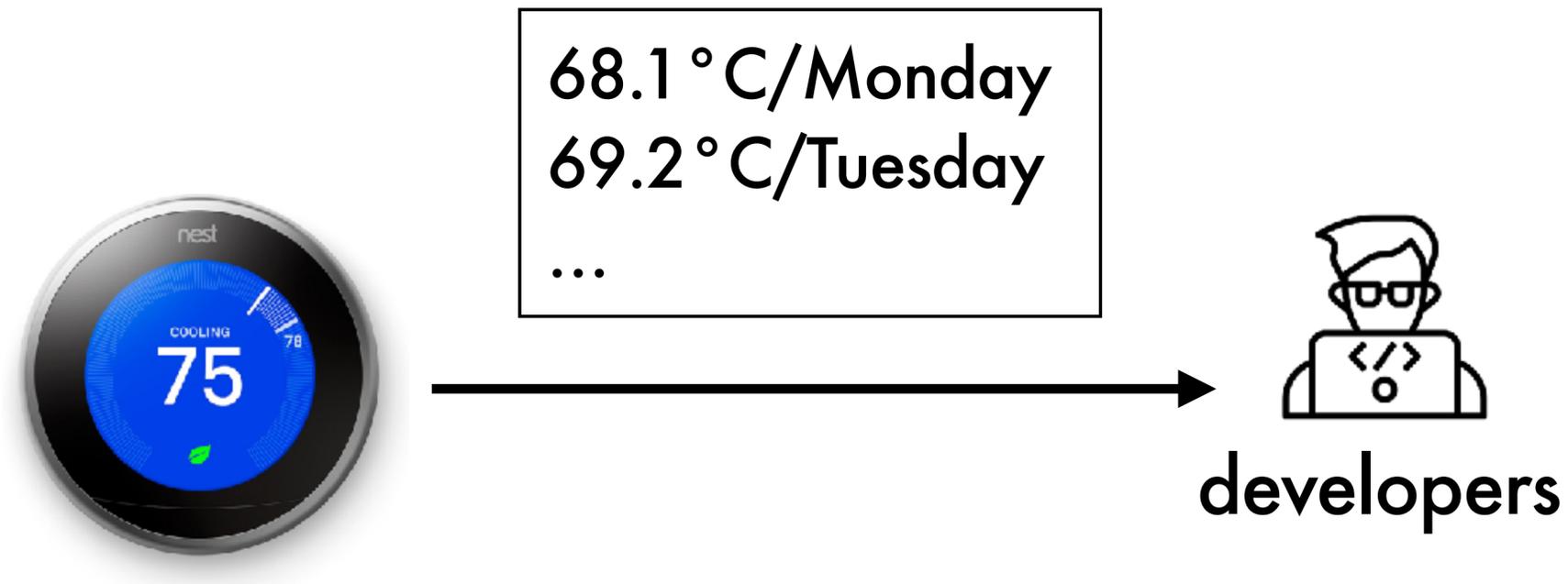
iOS 13



Talk outline

1. Modular Privacy Flows (MPF) in a Nutshell
2. Why MPF
3. **How** MPF
4. **When and when not** MPF
5. **Future Work**

How can Nest prove that they only collect aggregated data?



Open source?

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



↑

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

2. Implement - Peekaboo

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

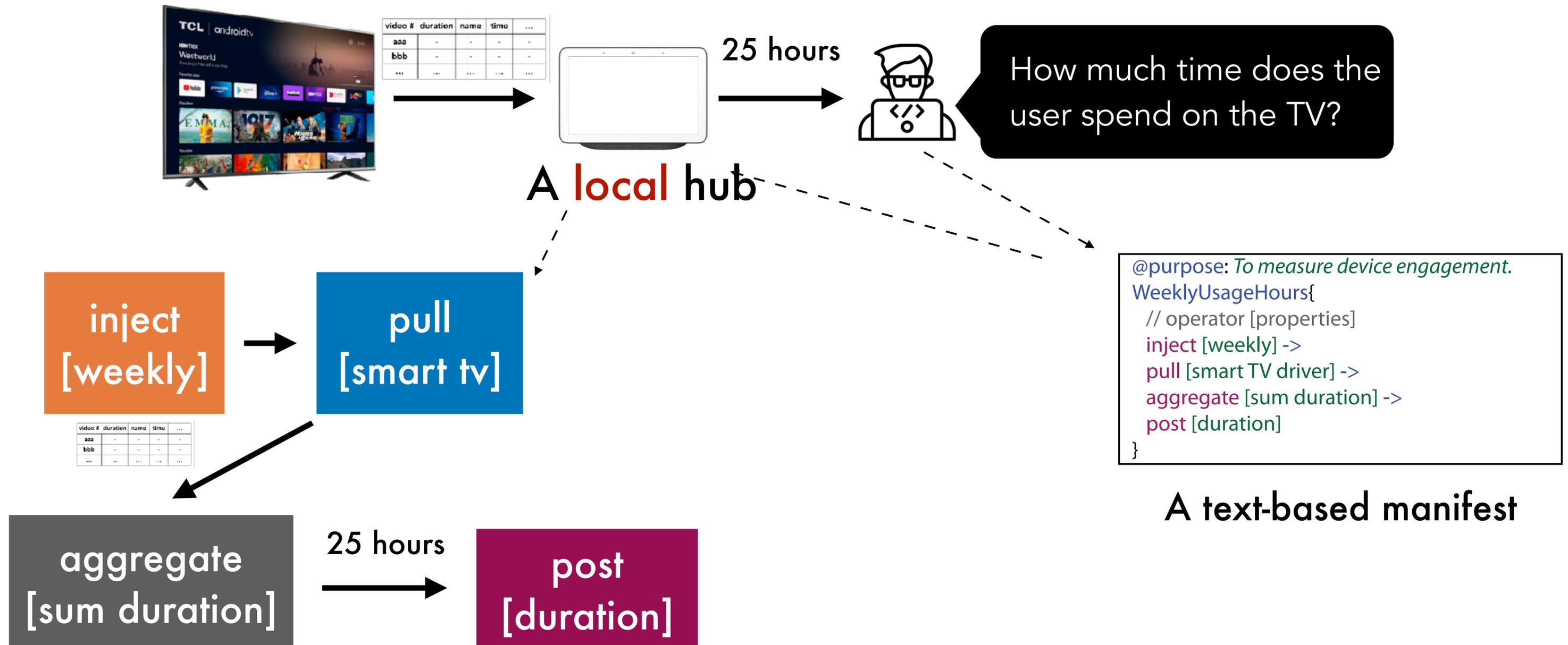


How much time does the user spend on the TV?

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

2. Implement - Peekaboo

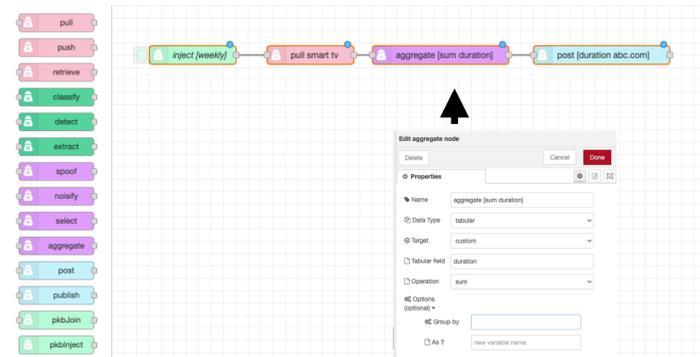
A trusted **runtime** with pre-loaded implementations



A trusted **runtime** with pre-loaded, open-source implementations



Smart home app store



Programming environment
with operators

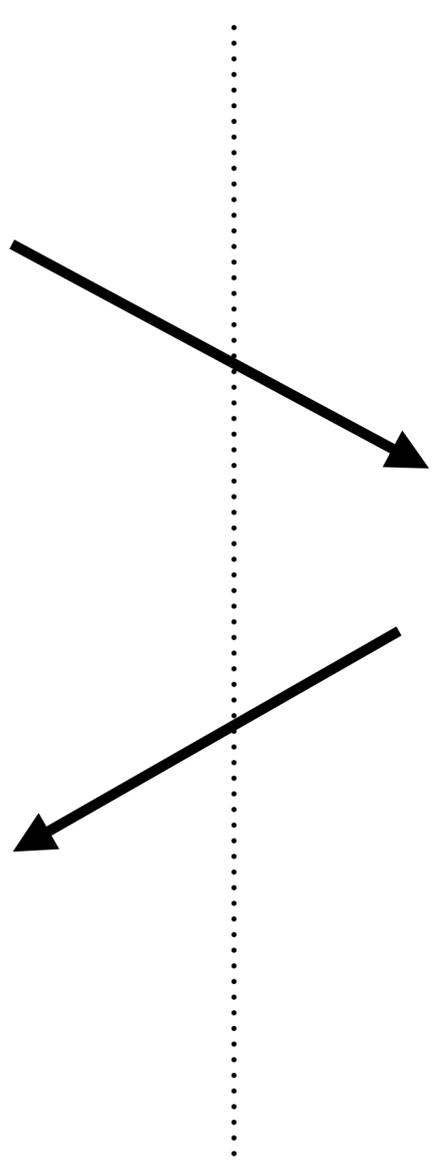


Runtime with preloaded
implementations

App developers

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

Manifest



Smart home hub → privacy firewall

Smart home app →

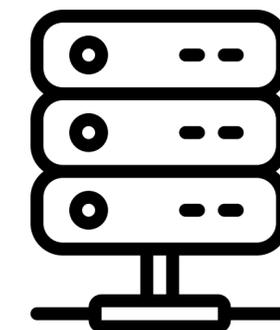
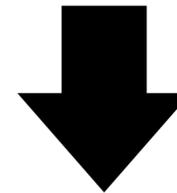
```
@purpose: To measure device engagement.  
WeeklyUsageHours{  
  // operator [properties]  
  inject [weekly] ->  
  pull [smart TV driver] ->  
  aggregate [sum duration] ->  
  post [duration]  
}
```



Edge devices



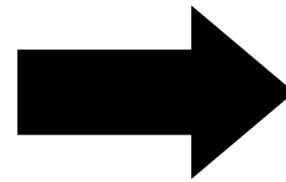
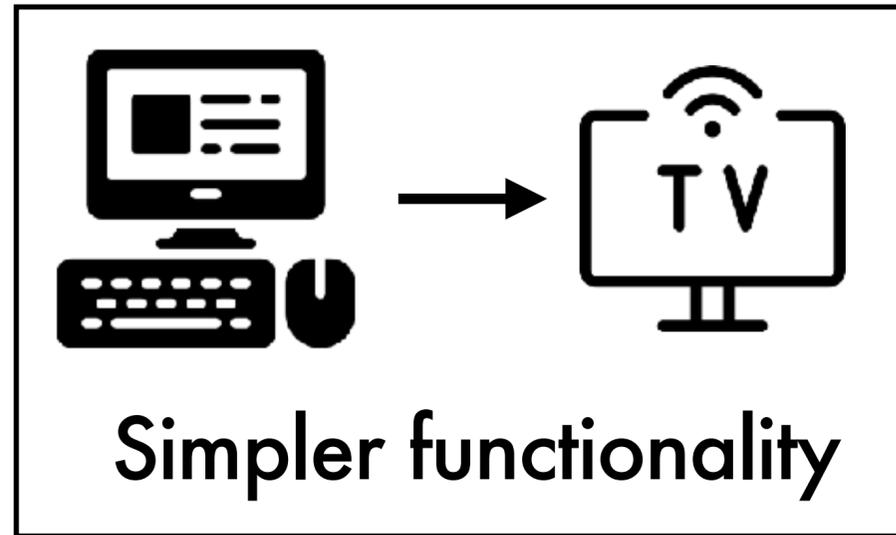
A local hub



Cloud

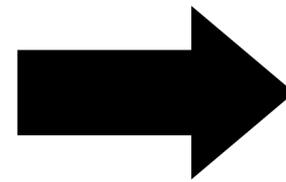
“Privacy firewall”

Peekaboo v.s. Firewall



Whitelist-only
Developer-in-the-loop

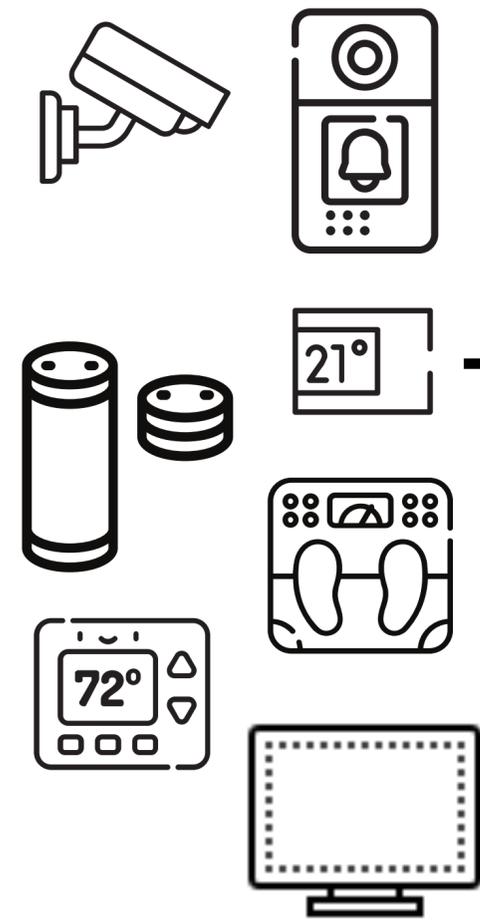
77% Apps do not
need raw data.



Pre-process users' data

How Peekaboo works

Handle heterogeneous hardware with device drivers



Device APIs

Device drivers

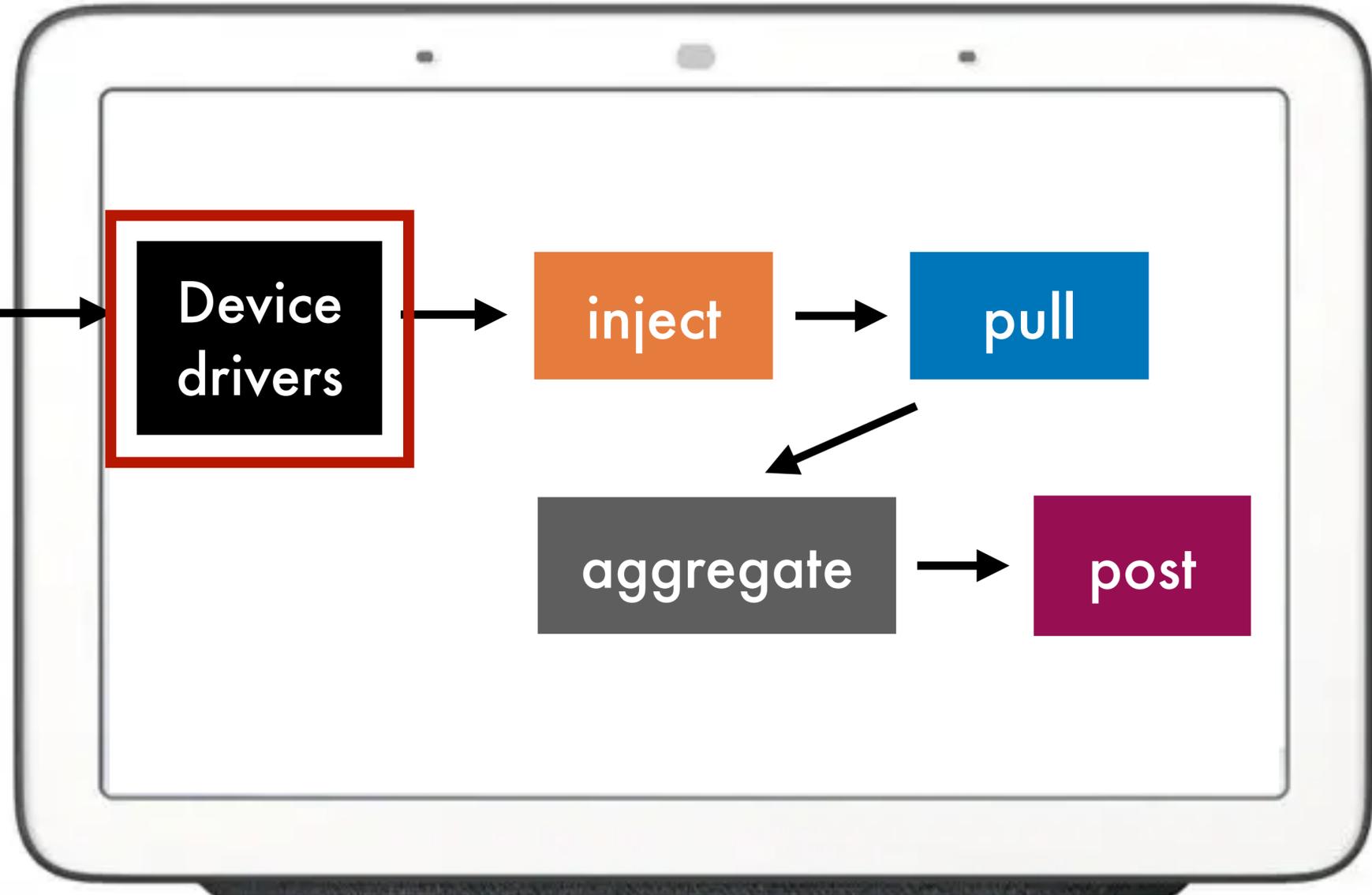
inject

pull

aggregate

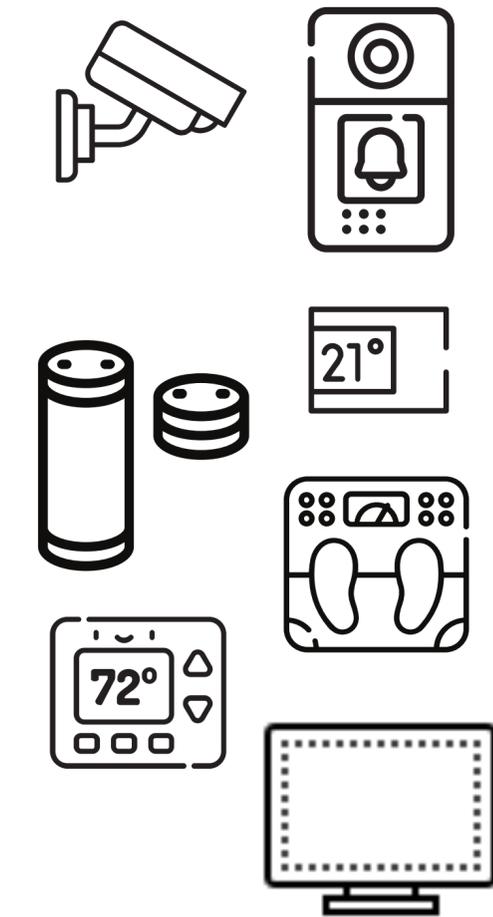
post

Edge devices



How Peekaboo works

A fixed set of operators



Edge devices

video, image, audio, tabular, scalar



A fixed set of operators

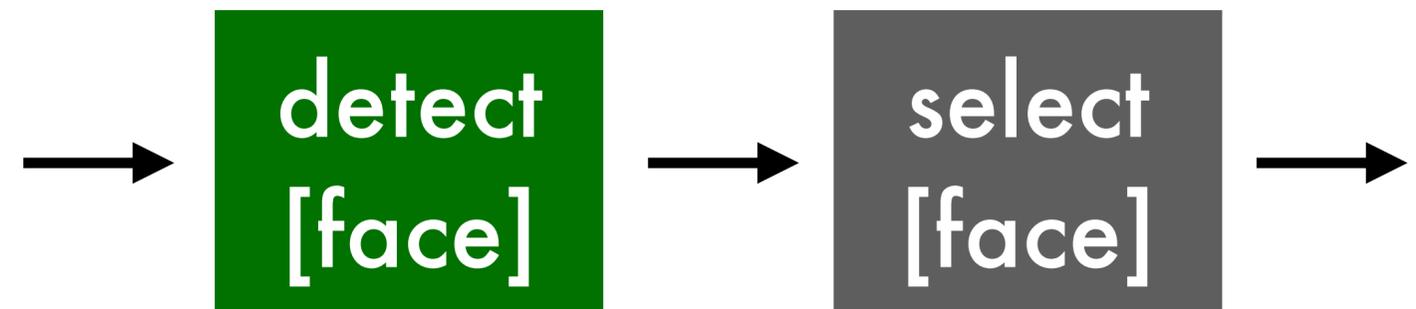
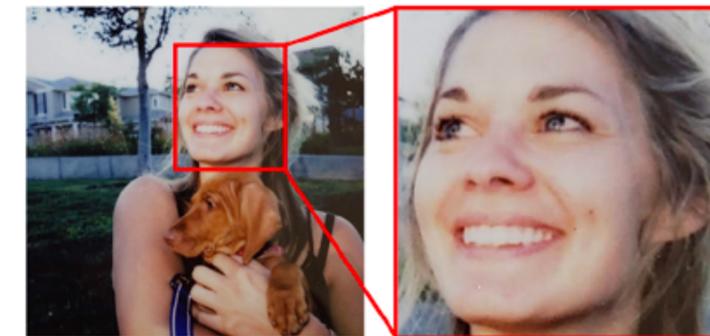
- pull
- push
- retrieve
- classify
- detect
- extract
- spoof
- noisify
- select
- aggregate
- post
- publish
- pkbJoin
- pkbInject

How Peekaboo works

An operator = A verb keyword

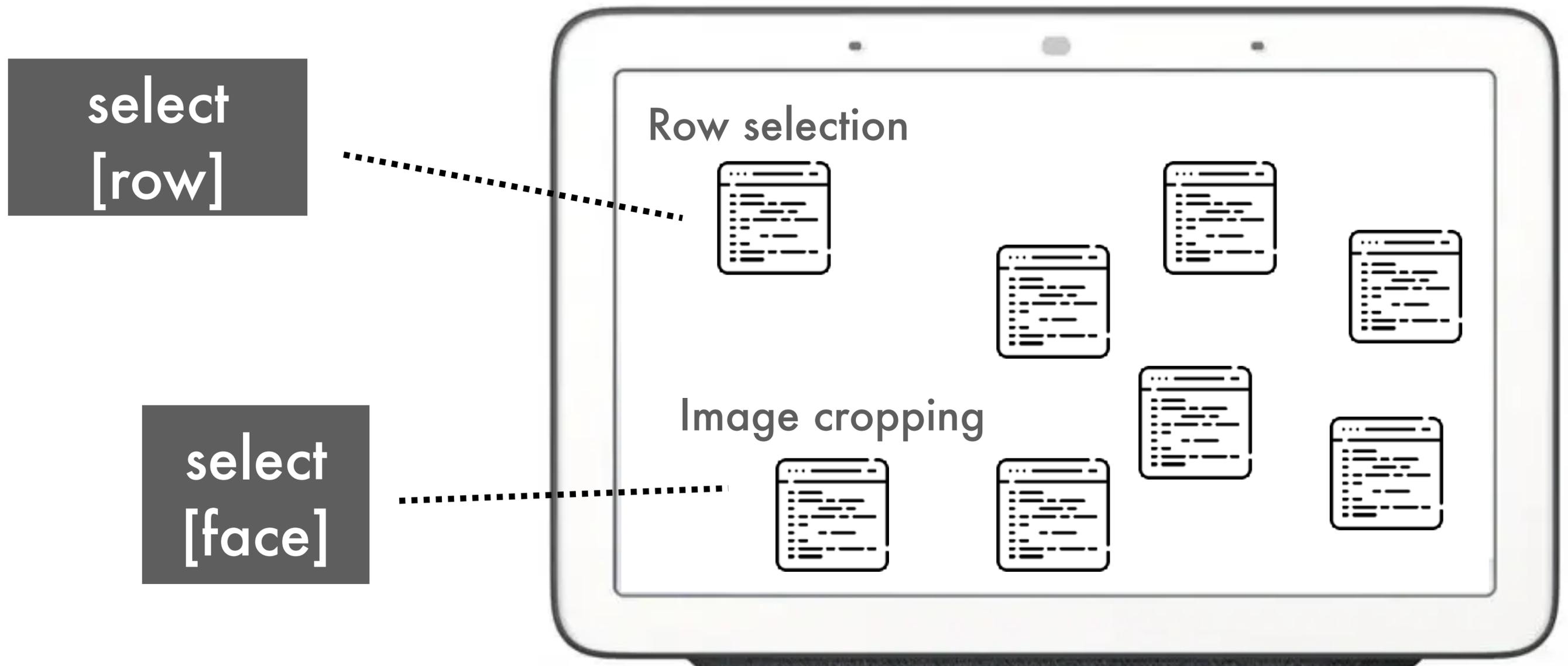
select
[row]

	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



How Peekaboo works

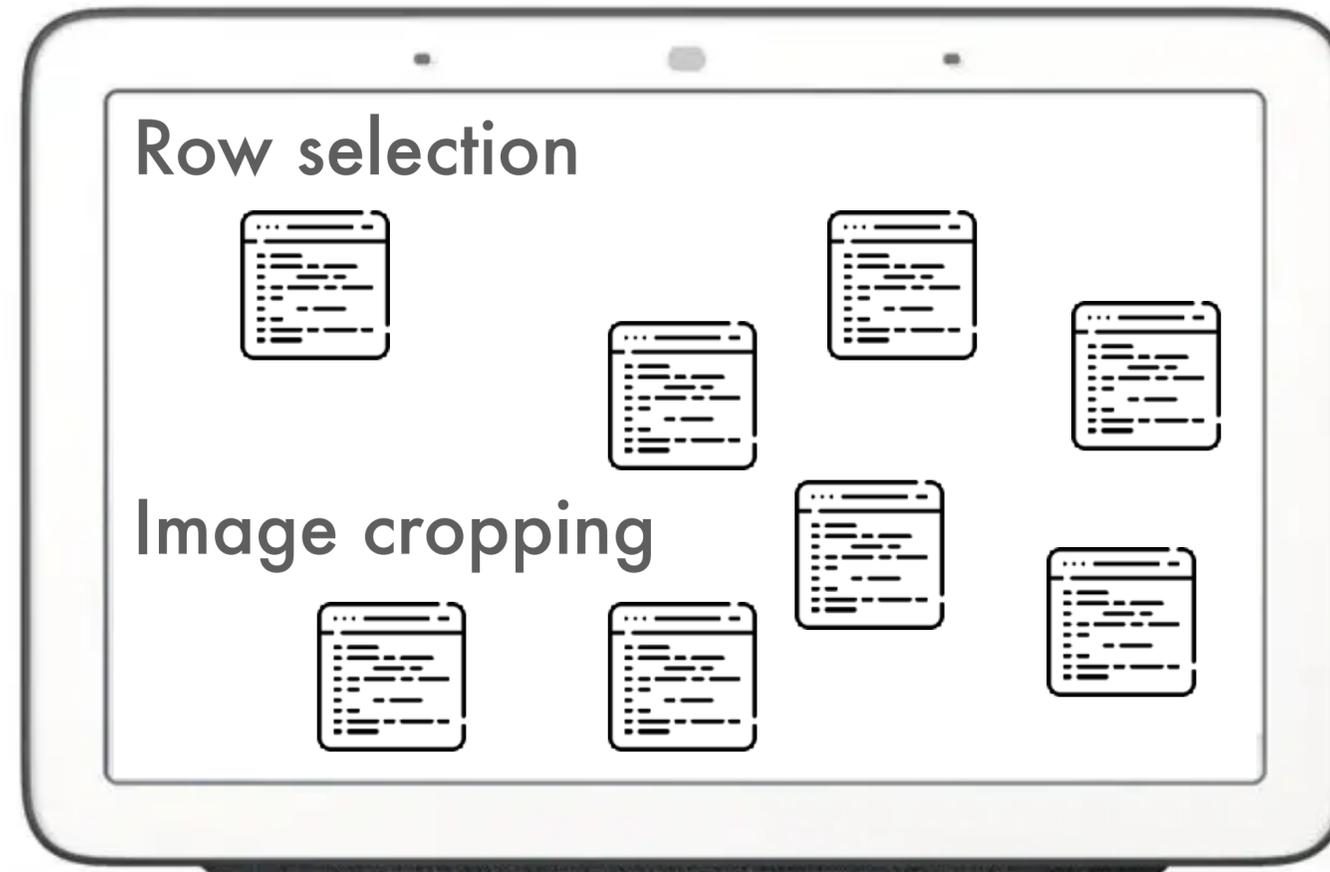
Operators are mapped to pre-loaded implementations



How Peekaboo works

A small set of pre-processing algorithms improve privacy

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



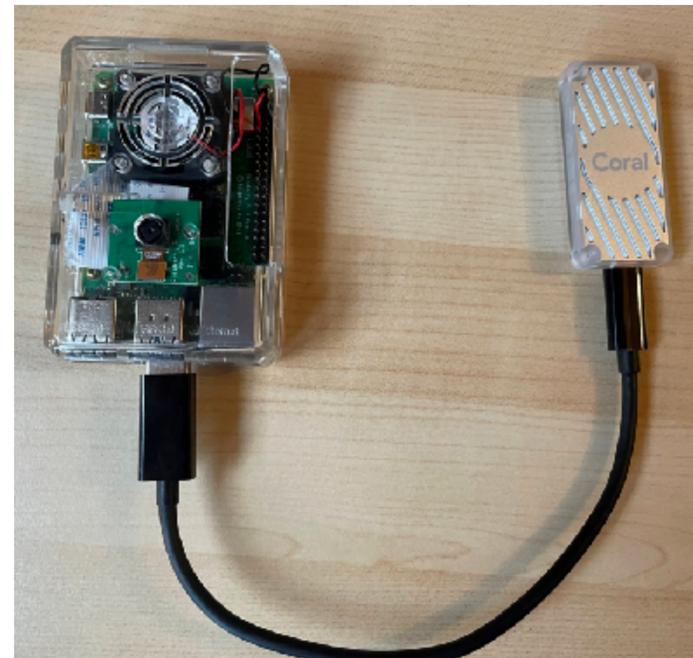
25 hours/week



Implementation (hardware)



Edge devices



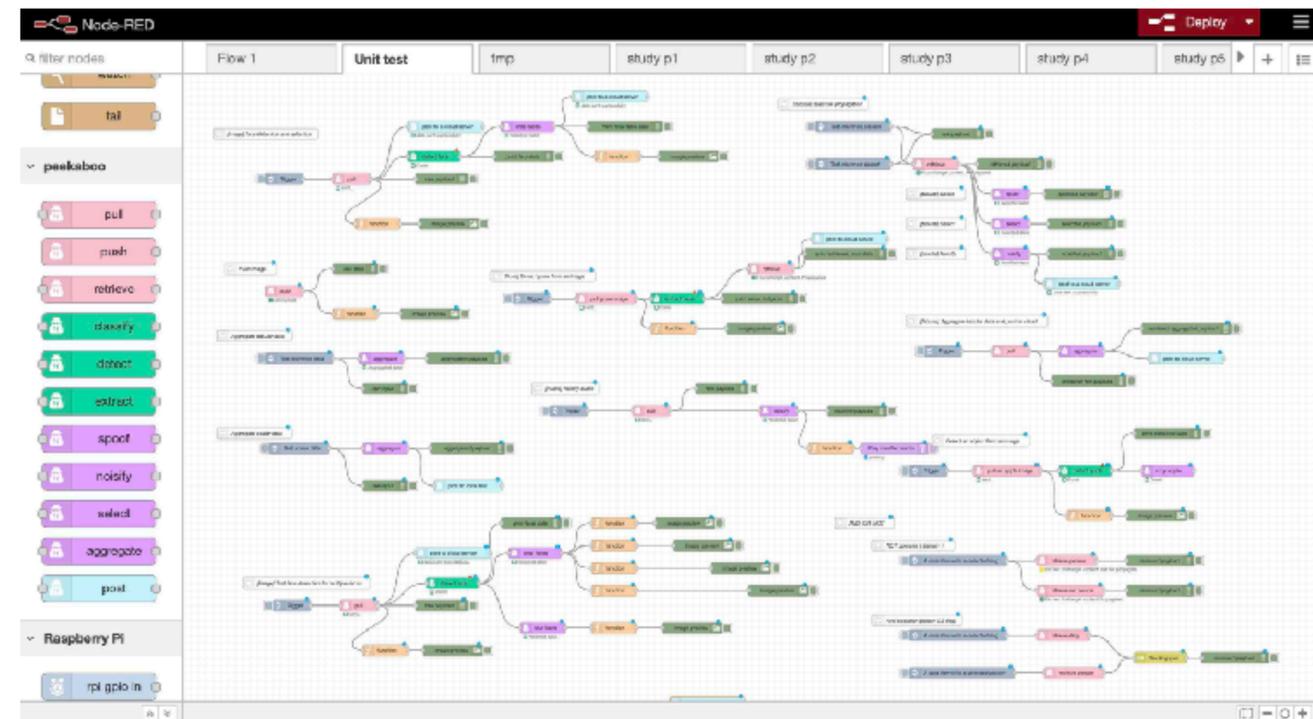
Raspberry Pi + TPU



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)

The screenshot displays the Node-RED web interface in a Chrome browser. The browser's address bar shows the URL `128.237.99.157:1880/#flow/82c51466.bda598`. The interface includes a sidebar with a search bar and a list of nodes under the category 'peekaboo'. The main workspace contains three distinct flows:

- Flow 1.1: water leak detection on the floor using cameras** - This flow is currently empty.
- Flow 1.2: water leak detection on the floor using microphones** - This flow consists of the following nodes in sequence: 'wait for pushed audio' (with a 'disconnect' indicator), 'recognize dripping sound', 'retrieve dripping event', 'blocking join (if appear 5 times)', 'pull audio', and 'post to cloud'.
- Flow 1.3: water leak detection on the floor using humidity sensor** - This flow consists of the following nodes in sequence: 'trigger every 30 mins', 'pull humidity', 'aggregate humidity scores', 'classify dripping sound', 'retrieve dripping event', and 'post to cloud'.
- Flow 2.1: Home inventory tracking using RFID** - This flow consists of the following nodes in sequence: 'trigger every 30 mins', 'pull RFID', 'select UUID', and 'post to cloud'.

The Node-RED sidebar on the left shows the following nodes under 'peekaboo': pull, push, retrieve, classify, detect, extract, spoof, noisify, select, aggregate, post, pkbJoin, and pkbinject.

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

Utility privacy tradeoff example



incognito voice assistant

6 speakers
112 audio files [1]

noisify

<5% random pitch shift



Microsoft
Cognitive Services

Speech word error rate:
9.27% → 11.88%

Speaker recognition:
100% → 27.7%

Developer studies



Task descriptions
IDE & Unit tests

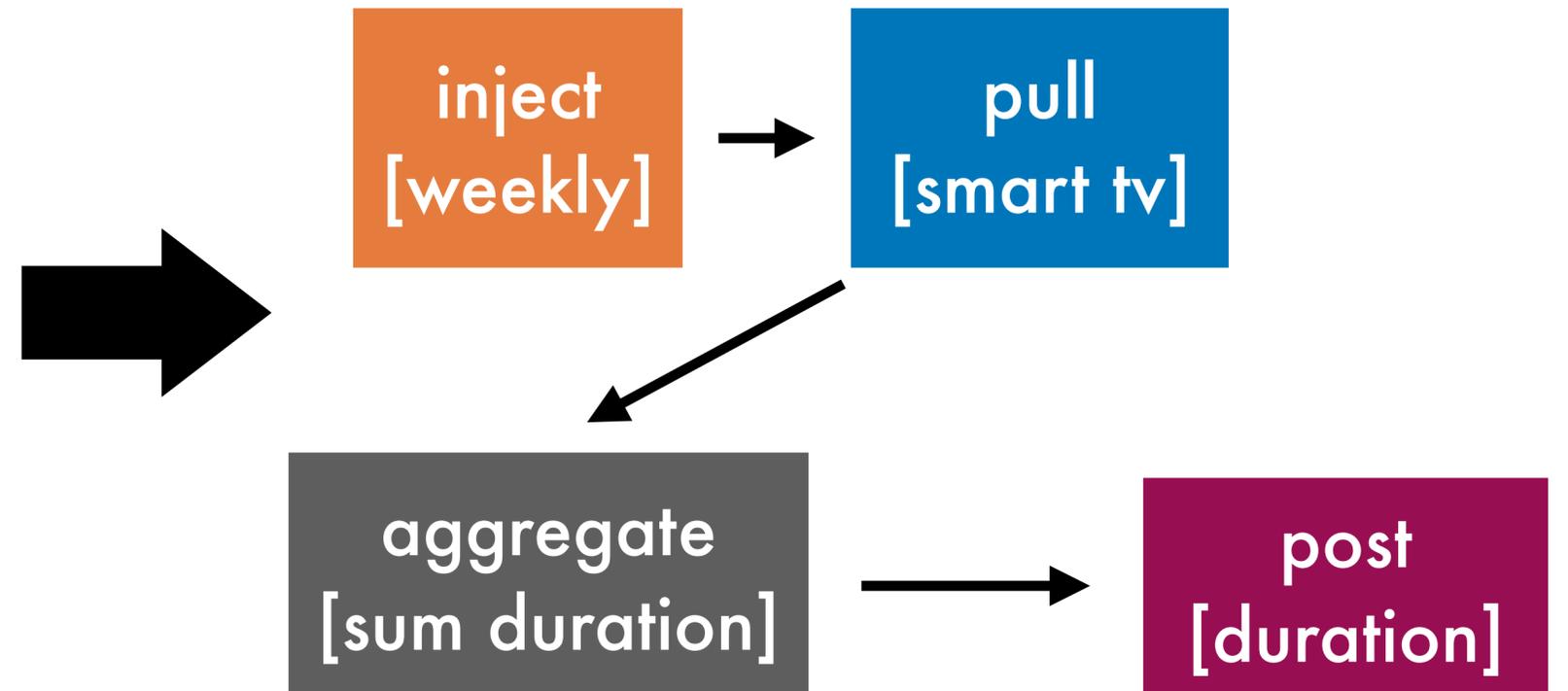
6 - 15 mins to
author a manifest

Advantages

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]
}
```

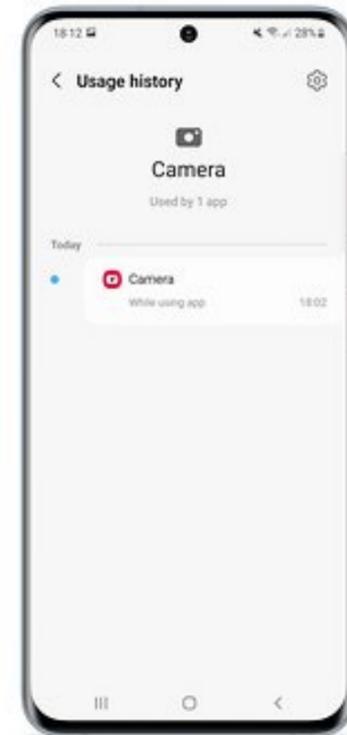
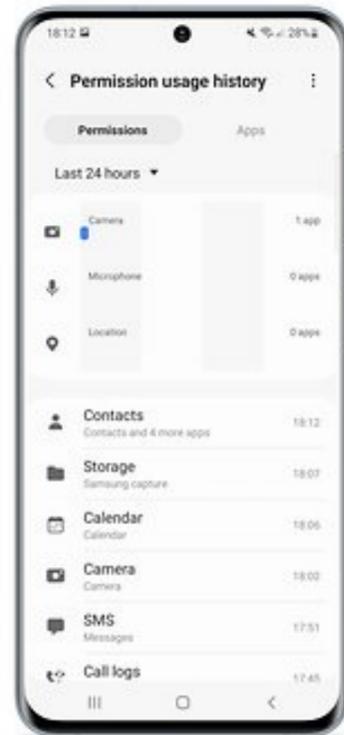
public, non-proprietary



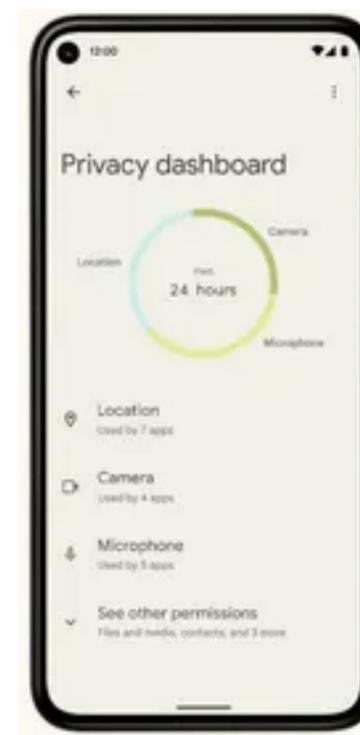
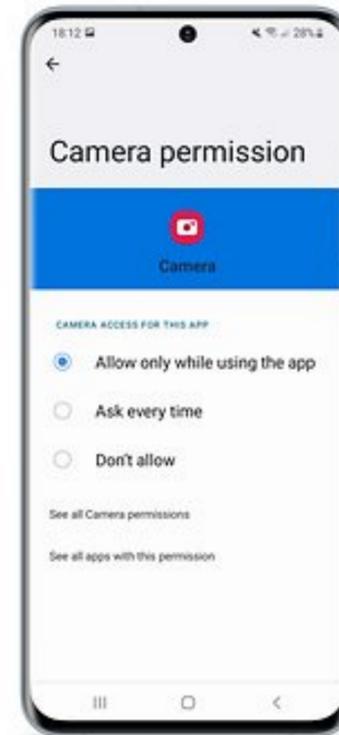
Advantages

Repetitive implementation and *distributed* interfaces

Samsung



Nest



Small developers?

Users?

Advantages

Manifests → *enforceable/dynamic* privacy nutrition labels

```
@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
<hr/>	
Total outgoing data packets	80
<hr/>	
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
<hr/>	
Customizations	
Rate limiting	N/A
More options

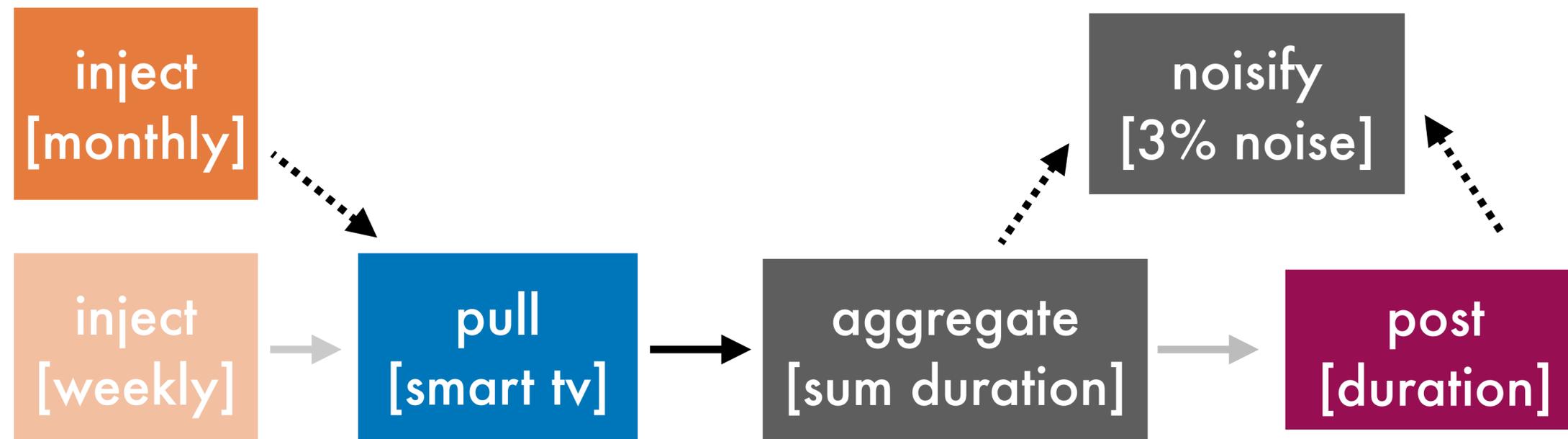
[1]

Advantages

Built-in fine-grained control through manifest rewriting

<u>Data Collection Disclosure</u>	
TV Usage Summary App	
<hr/>	
	<u>Customizations</u>
Rate limiting	N/A
More options

Change the rate
to **monthly**



Revisit: The permission granularity dilemma

More fine-grained permissions

→ Better privacy

→ More management burden for users

Harder learning curve for app developers

More implementation efforts for system builders

More coarse-grained permissions

→ Worse privacy

→ Overaccess risks

More users deny data requests

More complaints for system builders

Hard to gain trust from users for app developers

Revisit: The permission granularity dilemma

More fine-grained permissions.

→ Better default options.

Machine-readable permissions

→ Easier to audit.

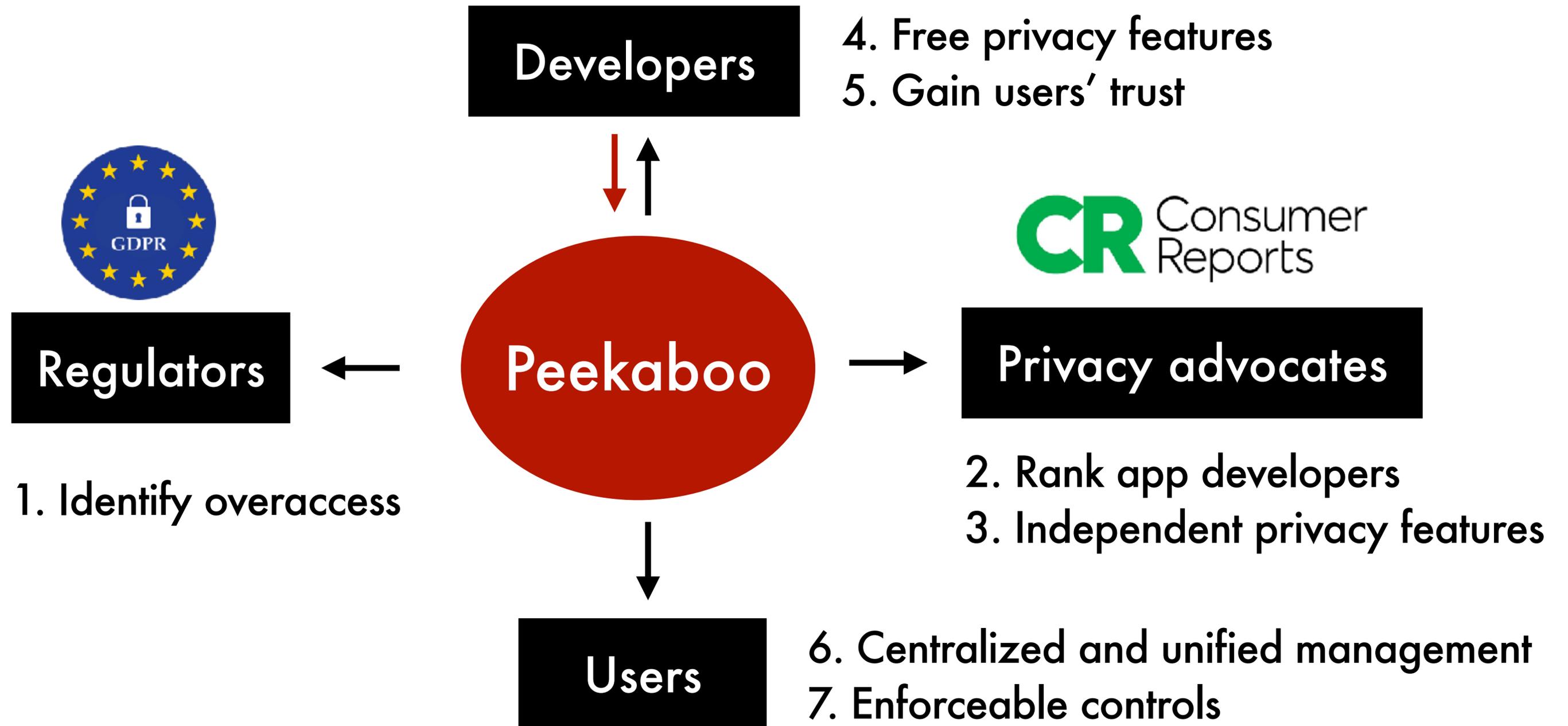
→ Better ecosystem. Good privacy drive-out bad privacy.

→ Aggregated management.

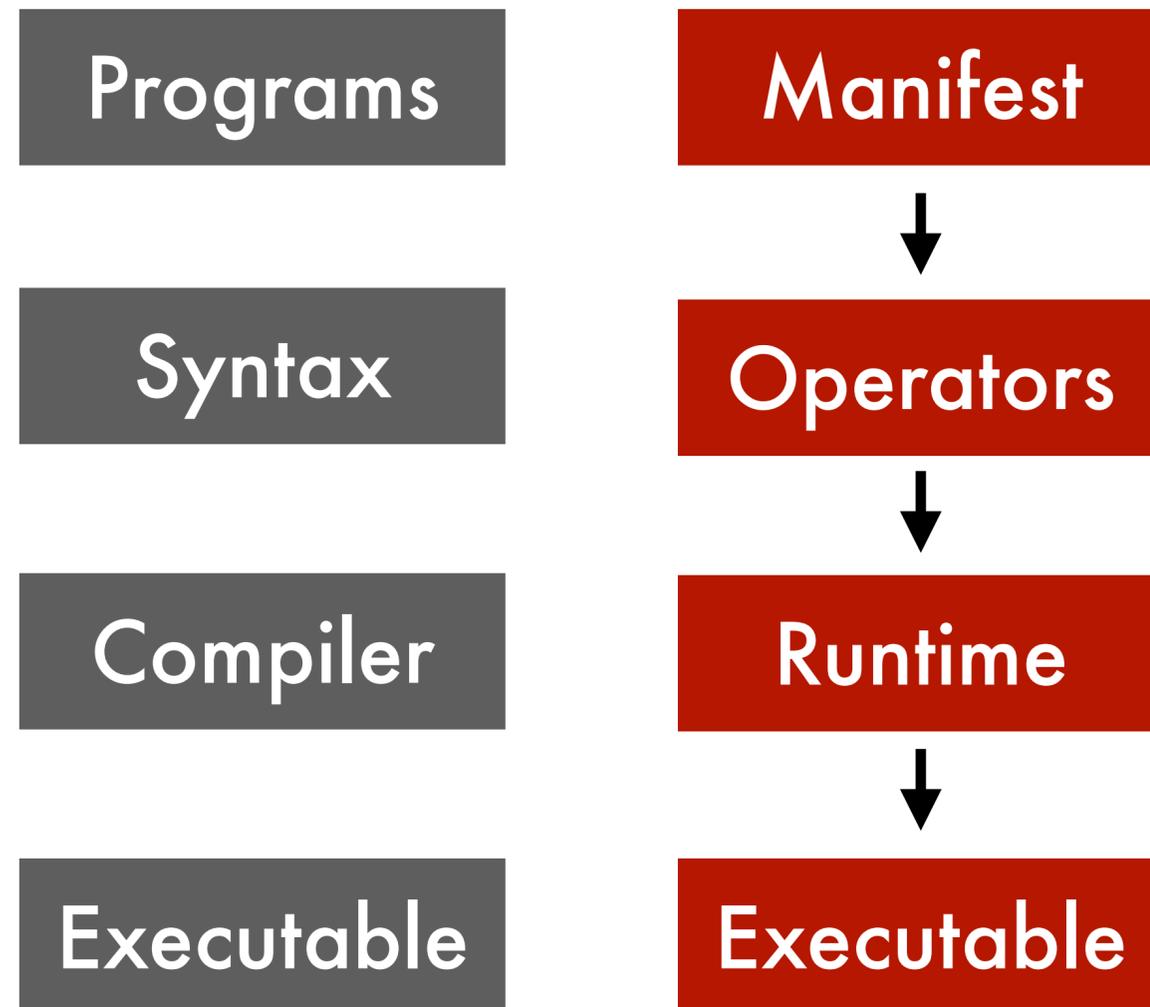
Decomposable (operator-based) permissions.

→ Fast development.

Let the good privacy drive out the bad privacy



MPF is a simpler compiler architecture.



A **fixed** set of operators

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

Talk outline

1. Modular Privacy Flows (MPF) in a Nutshell
2. Why MPF
3. How MPF
4. **When and when not MPF**
5. **Future Work**

Recap:

Privacy as *modular* information flow

Who (which app) sends the data?

Where the data is being sent to?

What data is being collected?

Why the data is being collected?

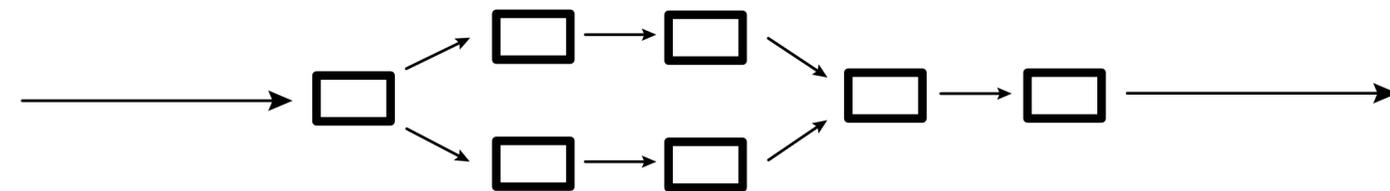
How the data is being stored?



Users



Developers



Flow-based programming

Future work: Broader application domains



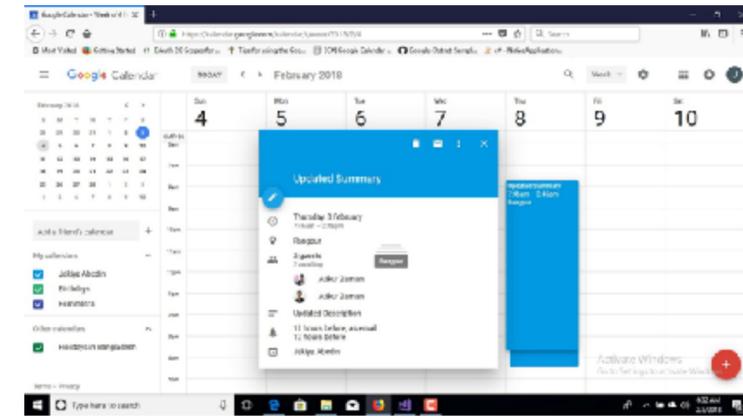
Peekaboo



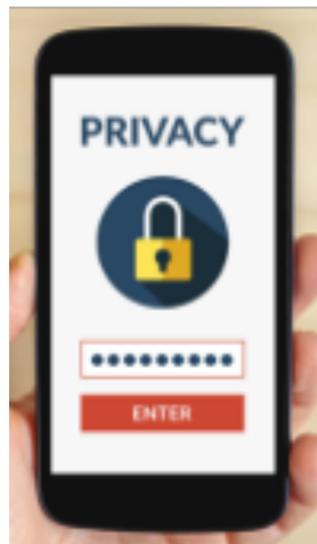
Smart City



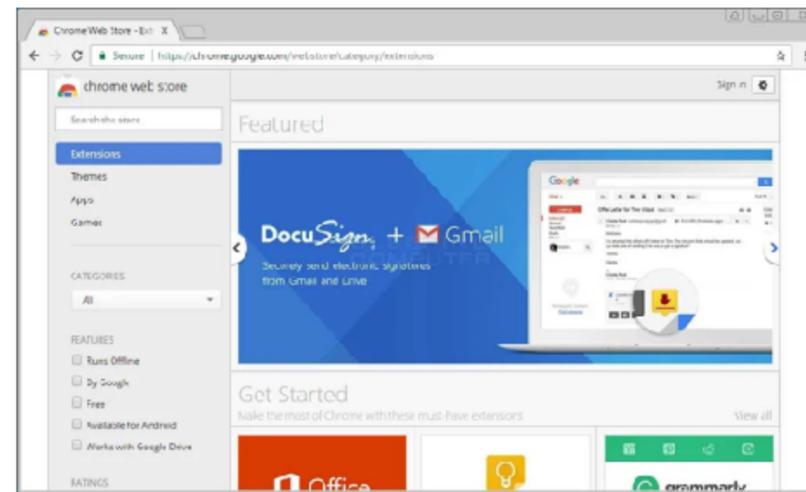
Social network?



Personal data API?



Mobile apps?



Chrome extensions?

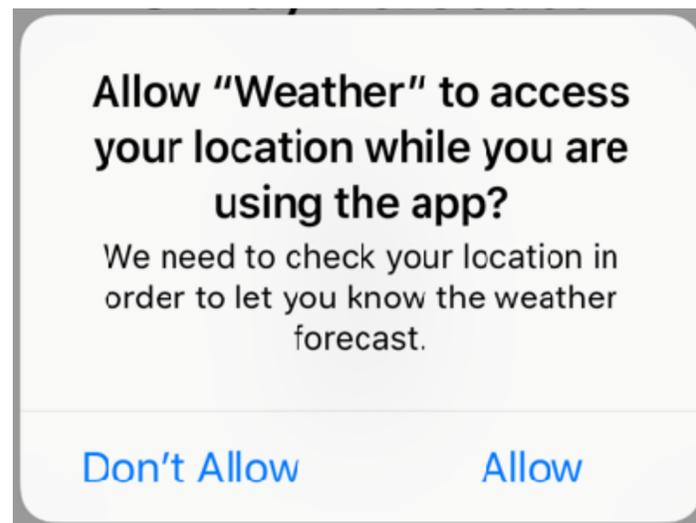


ChatGPT Plugin?

MPF v.s. Binary permissions

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

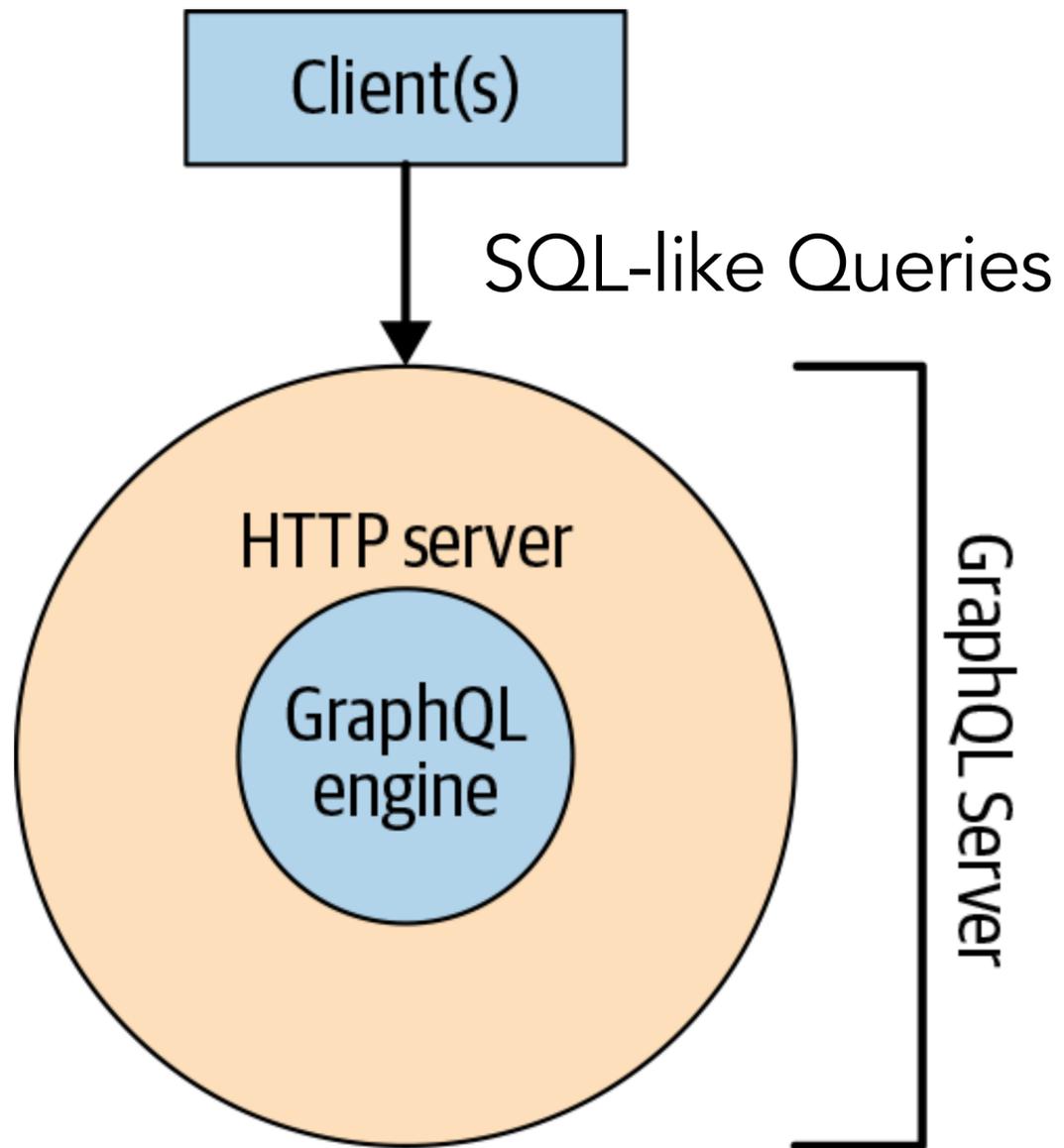
Android Permission Manifest



Popup window

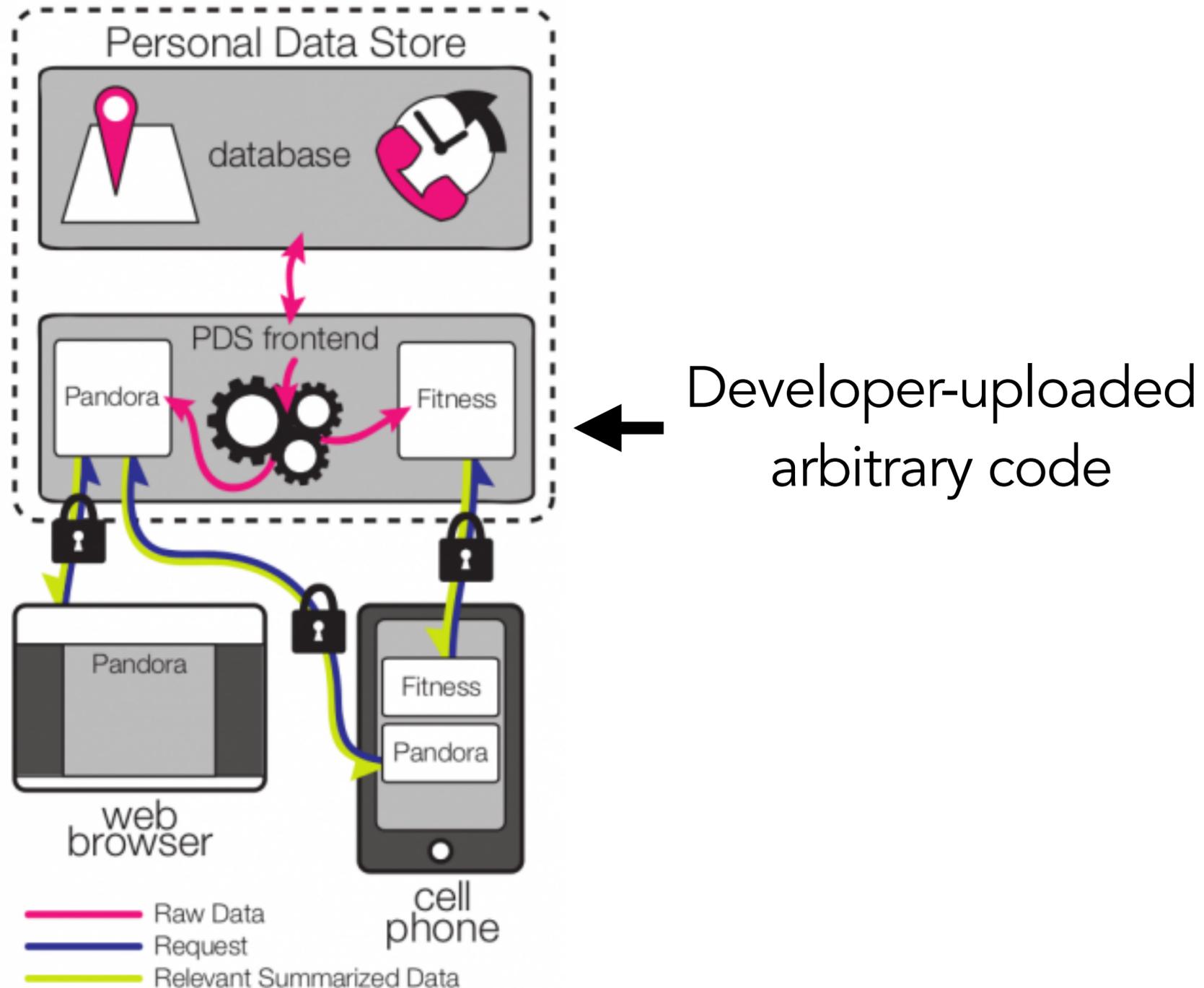
1. System implementation
2. API complexity
3. End-user management

MPF v.s. Database approaches (e.g., GraphQL)



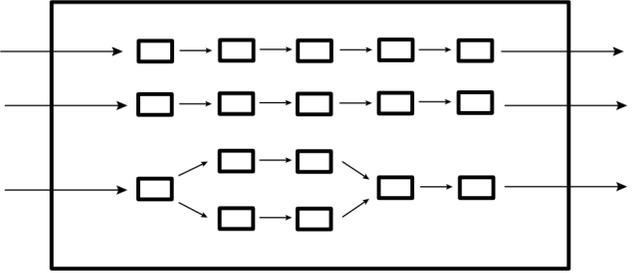
1. Flexibility/Extendability
2. Auditability

MPF v.s. Remote Code Execution

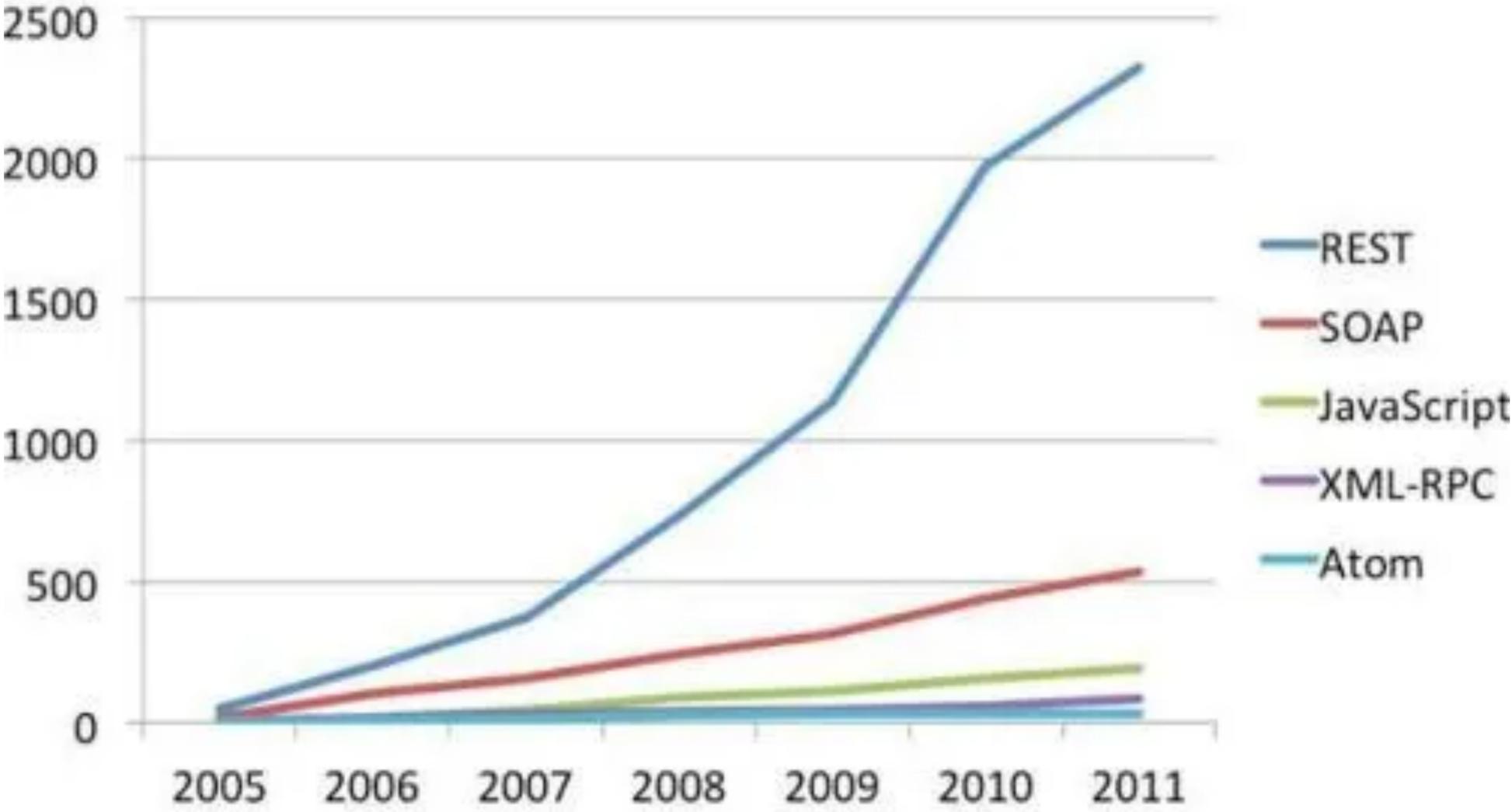


1. **Auditability**
2. **App development**
3. **Security**

Modular Privacy Flows today is REST in 2000.

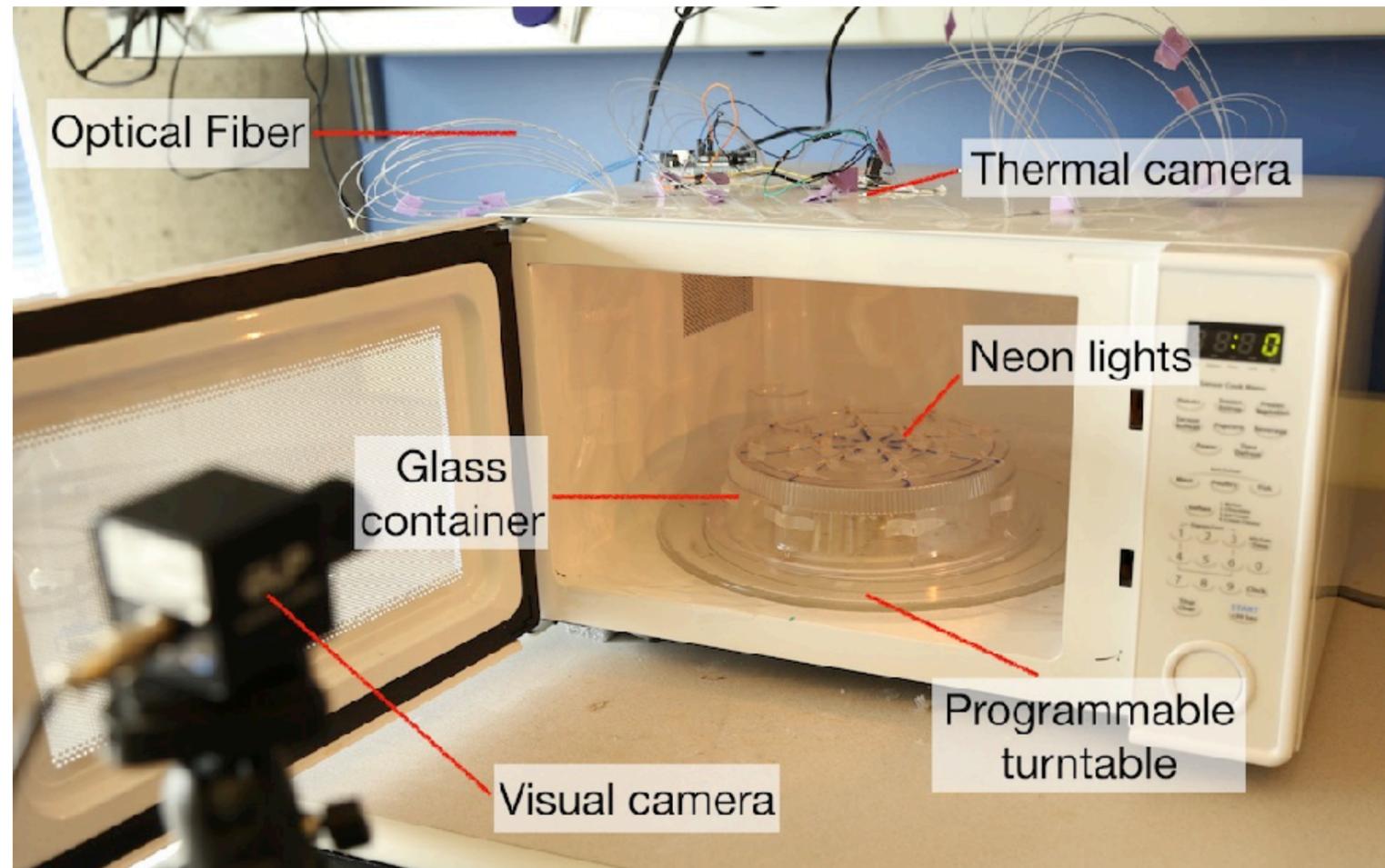


GET	<code>/pet/{petId}</code>	Find pet by ID
PUT	<code>/pet</code>	Update an existing pet
DELETE	<code>/pet/{petId}</code>	Deletes a pet
POST	<code>/pet/{petId}/uploadImage</code>	uploads an image



A story behind Modular Privacy Flows

Software-defined Cooking



No Turntable



Default Turntable

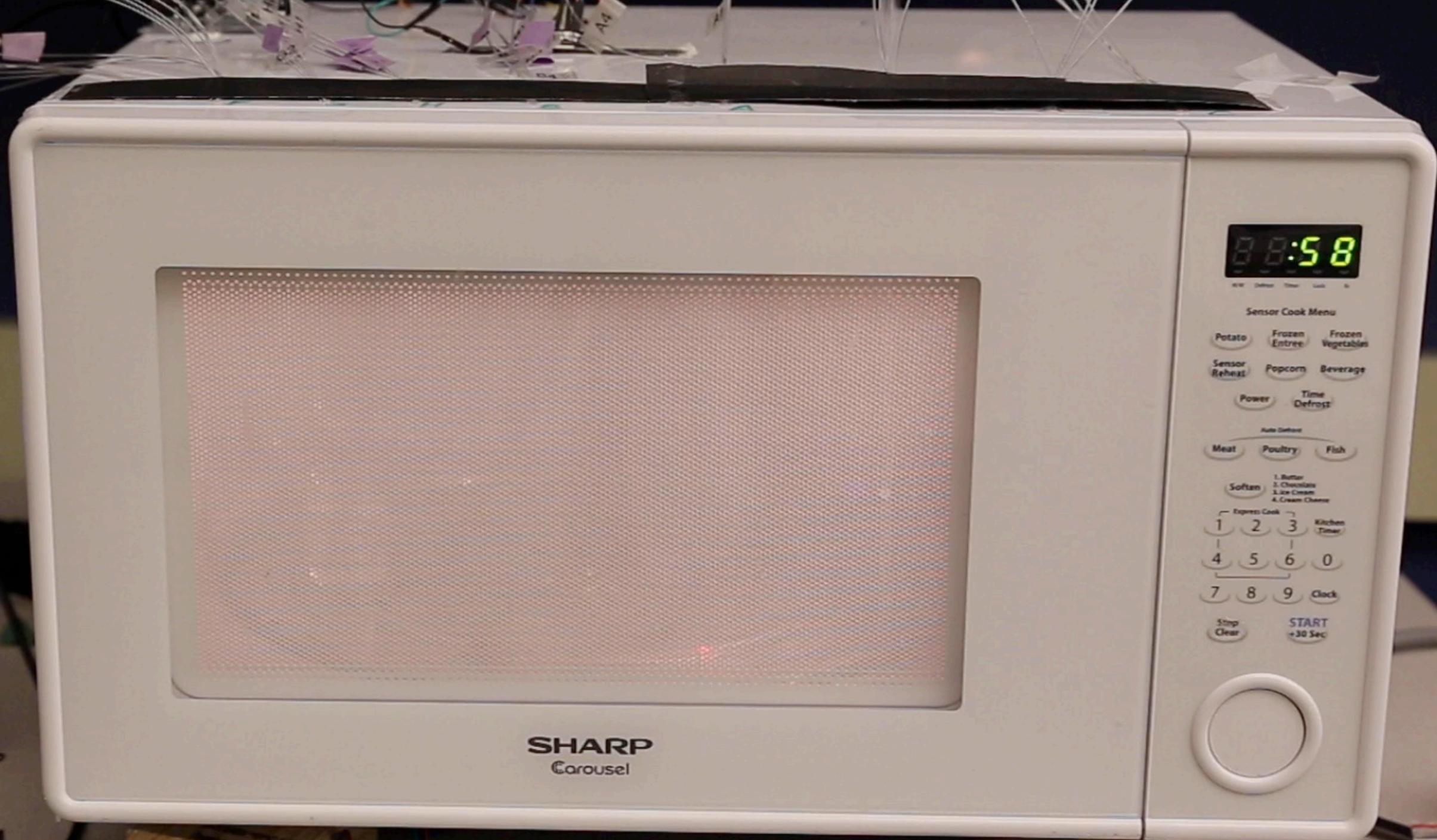


SDC Uniform Heating



SDC Arbitrary Heating





SHARP
Carousel

8:58

Sensor Cook Menu

Potato Frozen Entree Frozen Vegetables
Sensor Reheat Popcorn Beverage
Power Time Defrost

Auto Defrost

Meat Poultry Fish

Soften
1. Butter
2. Chocolate
3. Ice Cream
4. Cream Cheese

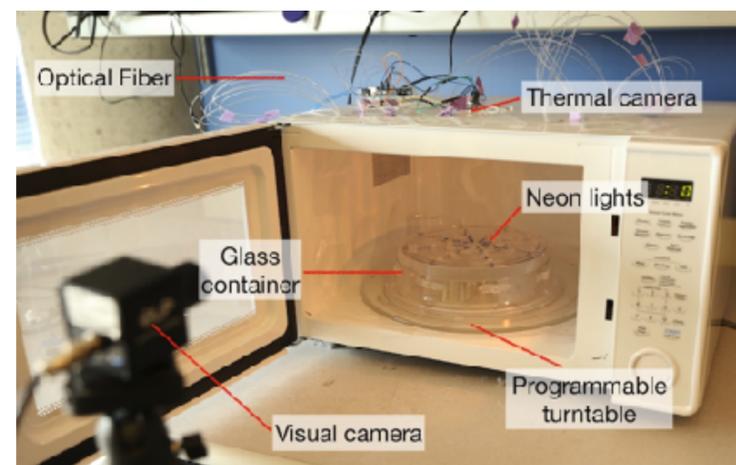
Express Cook

1 2 3 Kitchen Timer
4 5 6 0
7 8 9 Clock

Stop Clear START +30 Sec

A story behind Modular Privacy Flows

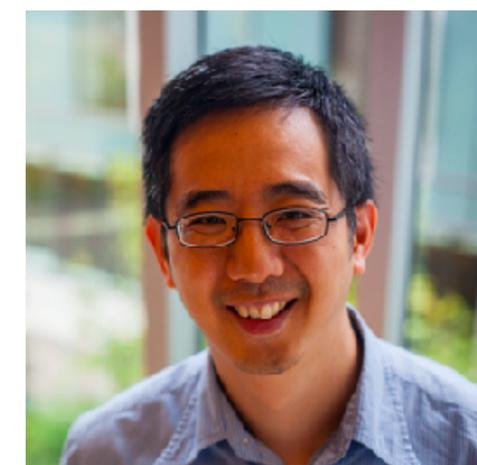
Software Defined
Hardware



+Operators



+Runtime



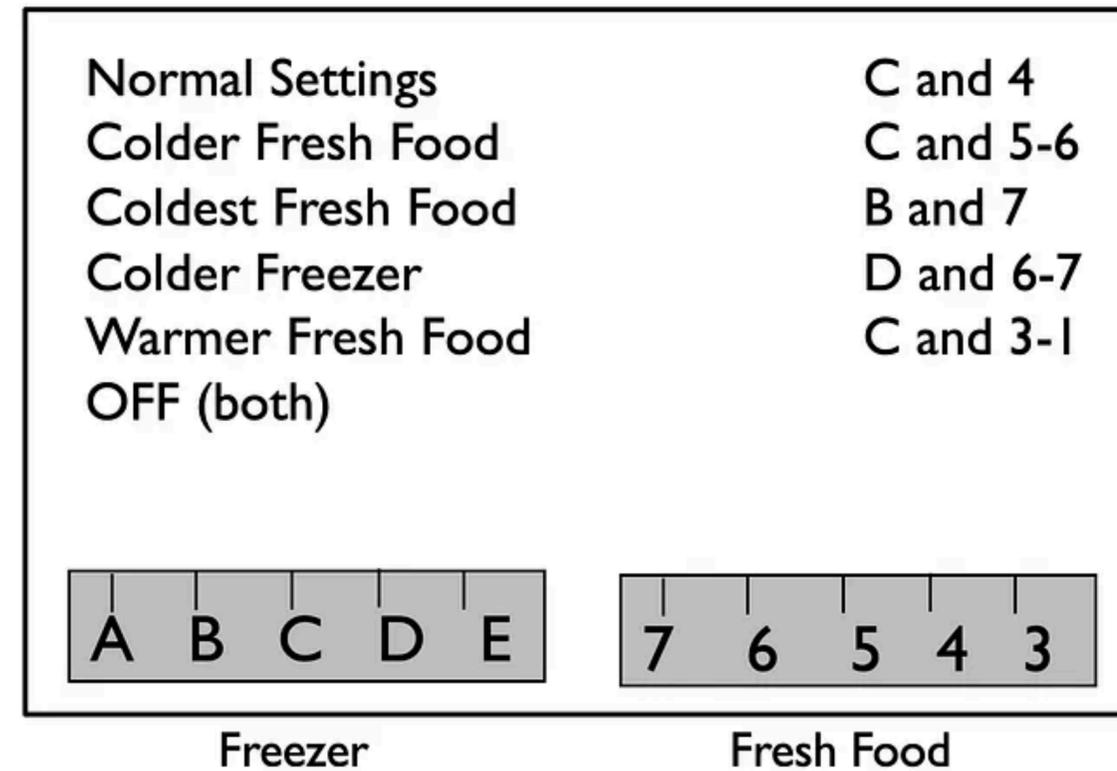
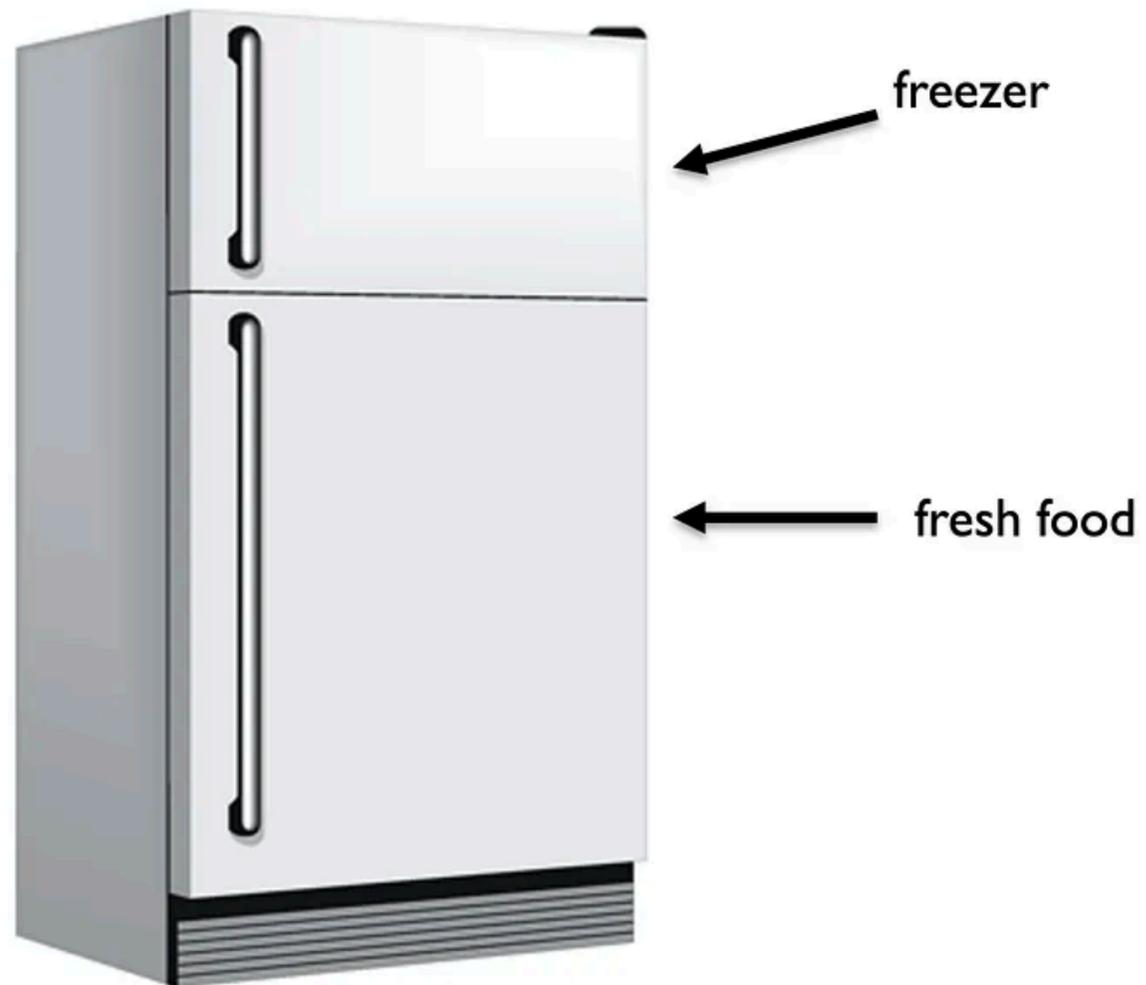
+Manifest

Talk outline

1. *Modular Privacy Flows (MPF) in a Nutshell*
2. *Why MPF*
3. *How MPF*
4. *When and when not MPF*
- 5. Future Work**

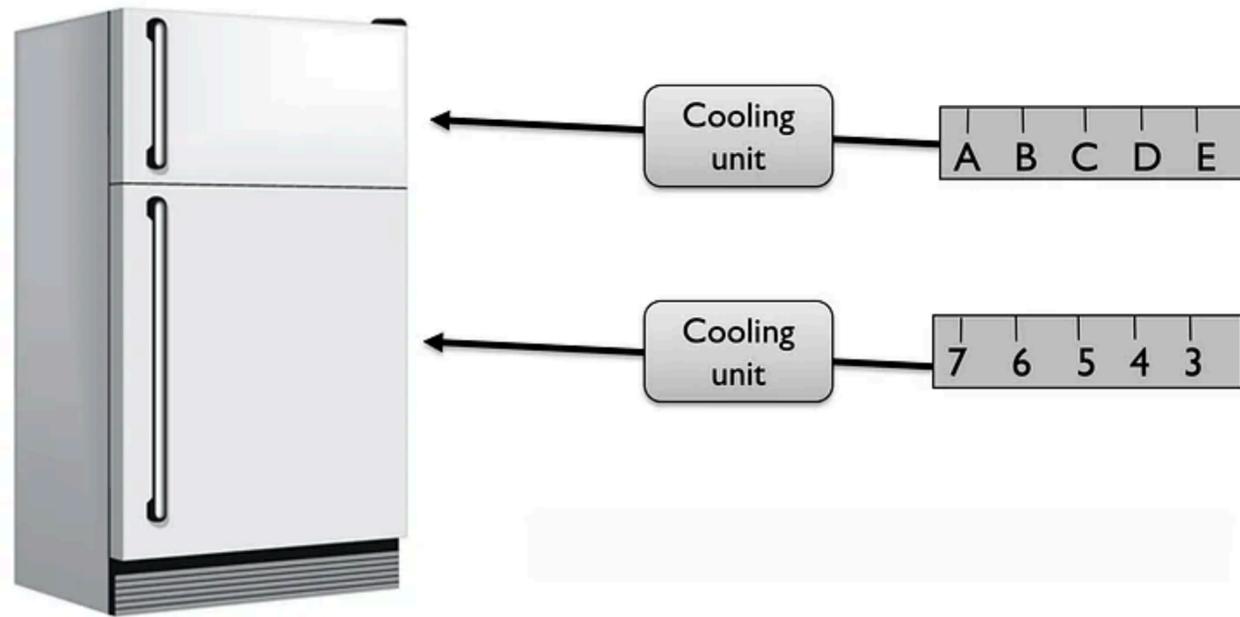
Computing systems are increasingly complex, we need something on the order of **a single page** demonstrating that the system will **work as intended**.

Mental models

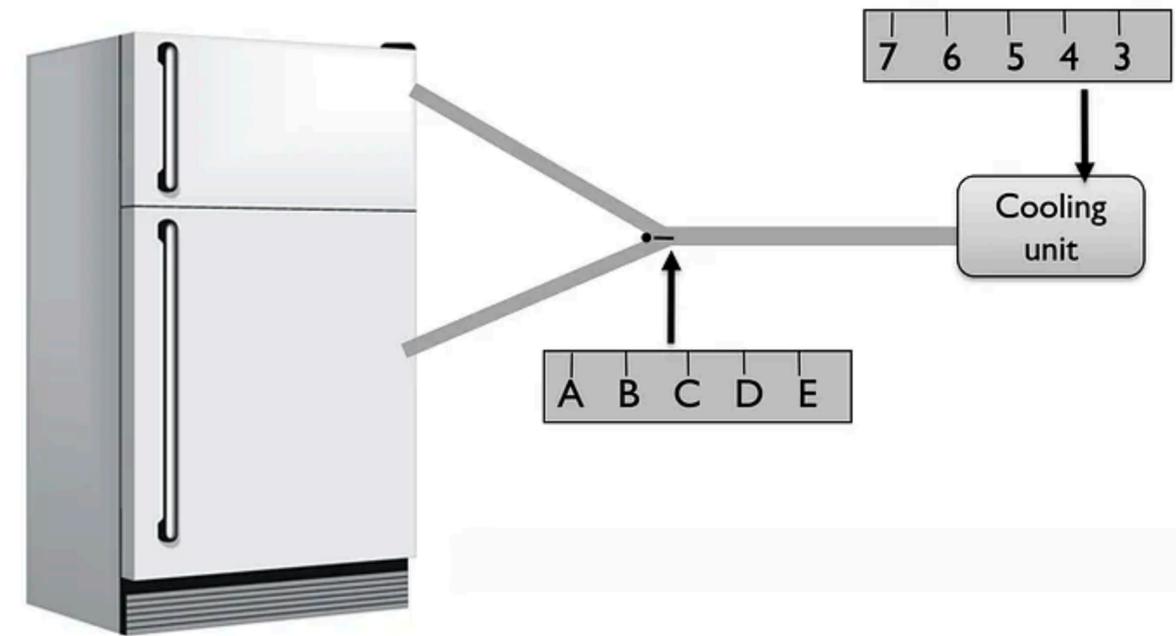


The Gap

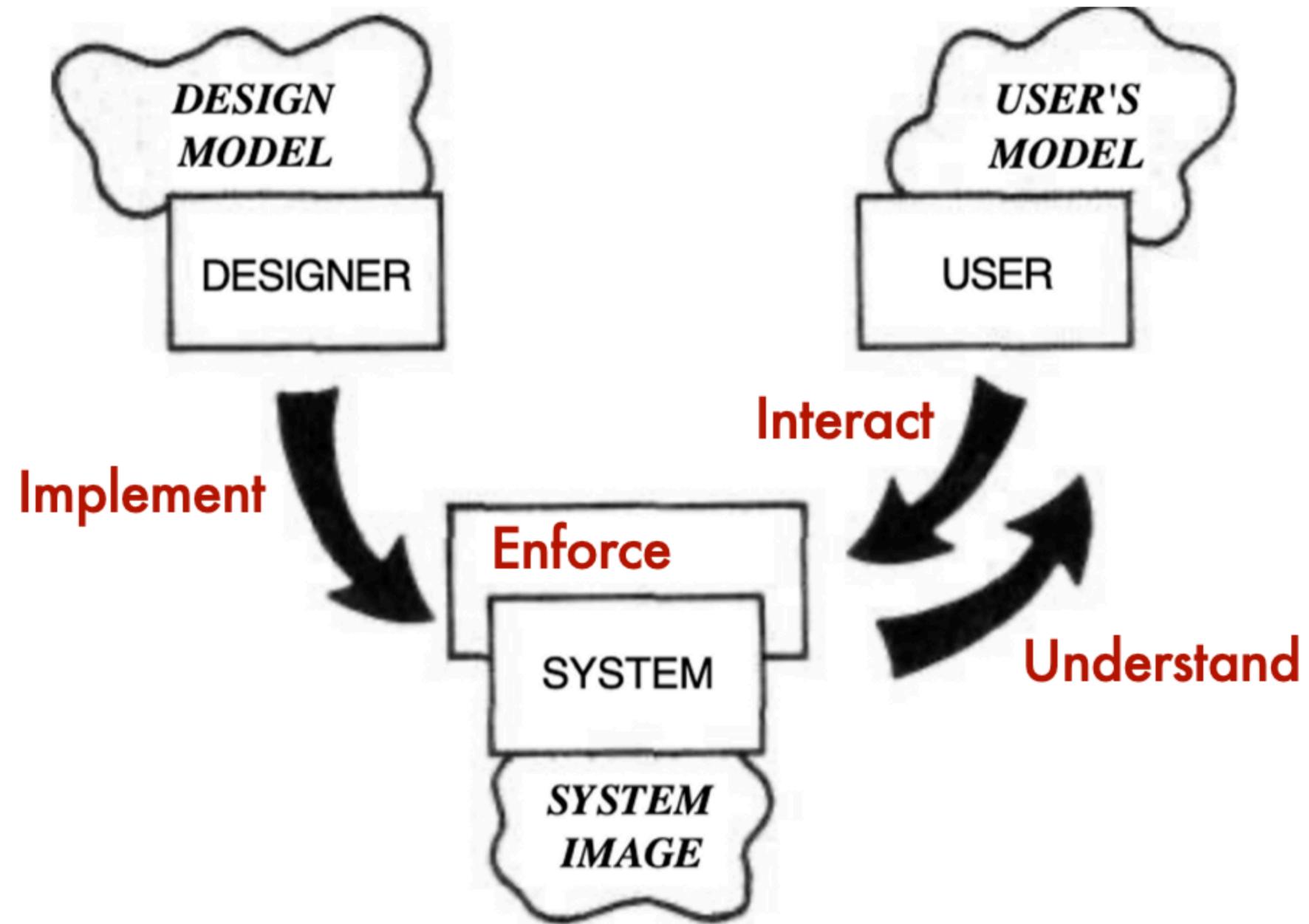
User's Conceptual Model



True System Model

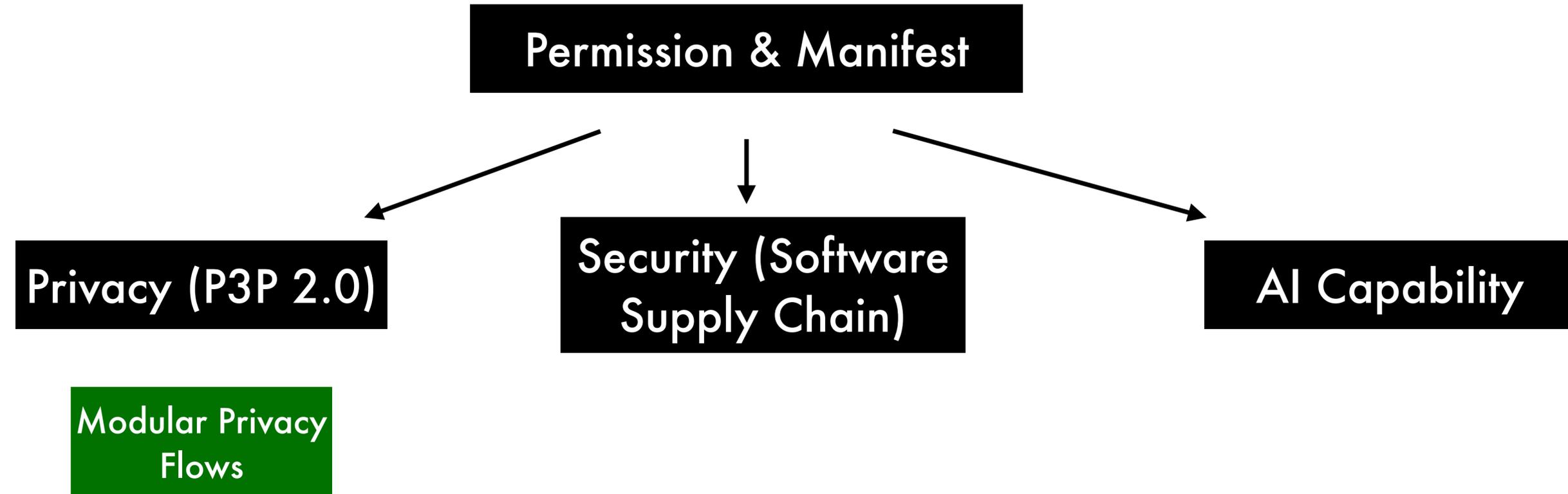


Our approach



Current research focus

Next generation permission system



Data Smith Lab is recruiting!

We study the **security and privacy of data systems** by researching the people who **design, implement, and use** these systems.

Contact: haojian@ucsd.edu
<http://haojianj.in/>

Acknowledgment

