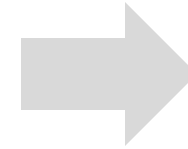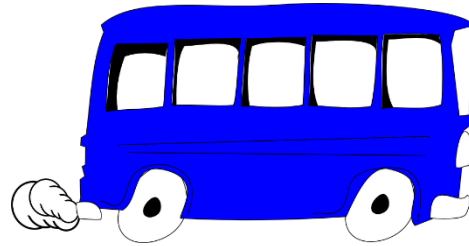# Data Infrastructure:
## Use Cases and Architecture

Joy Bonaguro

Chief Data Officer

City and County of San Francisco
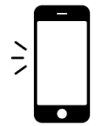
COIT Subcommittee, Feb 2, 2018

# Agenda

- Data infrastructure use cases
- Architectural choices (within and between)
- Benefits of a strategy

# DI Use Case: Move and process data between systems

City database or application

City database or application



Mobile device
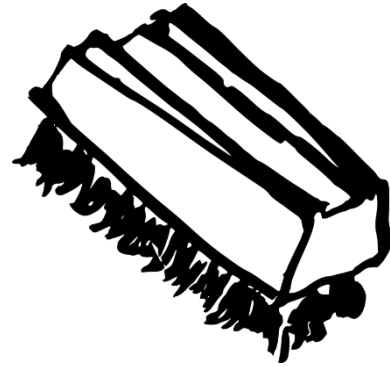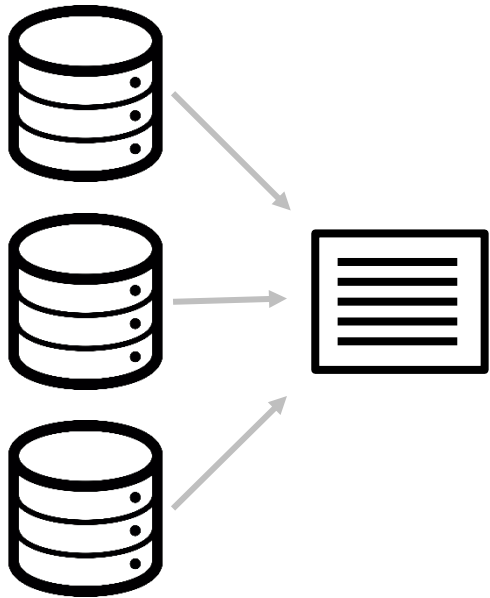
GPS recorder

Anything with a network connection

Data "Bus" processes and moves data around

# DI Use Case: Store data for use

City databases and applications
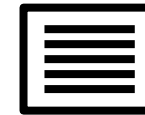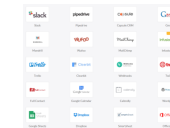


Data warehouse

Data lake

OD portal

TIMESCALE Specialized stores

API

Table views

FTP drop

BI tools

Web connectors

Data cube

**S1** Connect and extract data from source systems
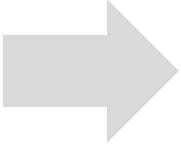
**S2** Clean and prepare data

**S3** Store data

**S4** Provide access

# DI Use Case: Monitor data quality
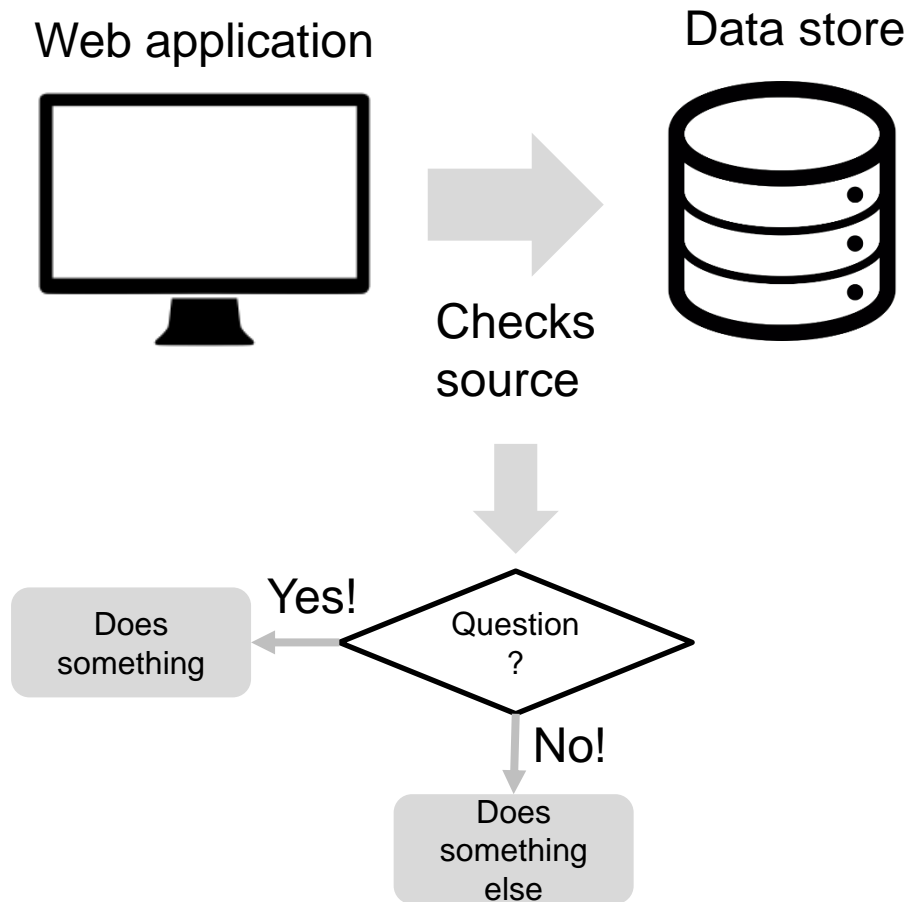
City data store

Database

Profiling…

Results Browser

Job: US Customer Data Profiling

| Input Field | Total Number | Minimum Length | Maximum Length | Minimum Value | Maximum Value |
|---|---|---|---|---|---|
| ID | 5438 | 9 | 9 | AAC434152 | ZZZ642455 |
| Name | 5438 | 11 | 39 | Anne Mullen | de Chana, Sergio Marques |
| Street | 5438 | 2 | 41 | # 3 Riverdrive Rd. East | Wilson & Kirk Road |
| City | 5438 | 3 | 20 | ABERDEEN | waterloo |
| State | 5438 | 2 | 2 | AB | WY |
| ZIP | 5438 | 4 | 10 | 01801-6202 | n2j4a9 |
| Country | 5438 | 1 | 13 | | United States |
| Phone | 5438 | 1 | 25 | (113) 072 3578 | x |
| Cell | 5438 | 4 | 14 | (113) 575 3765 | 9978 158 |
| Work | 5438 | 4 | 28 | (113) 007 6029 | x7562 |
| eMail | 5438 | 16 | 35 | Aaron.A.Koontz@thu.com | zoi.gibso@snomail.com |
| DoB | 5438 | 19 | 19 | Jan 1, 1900 12:00:00 AM | Mar 29, 2007 12:00:00 AM |
| Gender | 5438 | 1 | 1 | F | U |
| Active | 5438 | 1 | 1 | 0 | Y |
| CreditLimit | 5438 | 1 | 5 | 0 | 32800 |
| StartDate | 5438 | 19 | 19 | Apr 1, 2006 12:00:00 AM | Apr 1, 2009 12:00:00 AM |
| EndDate | 5438 | 19 | 19 | Apr 1, 2008 12:00:00 AM | Apr 1, 2014 12:00:00 AM |

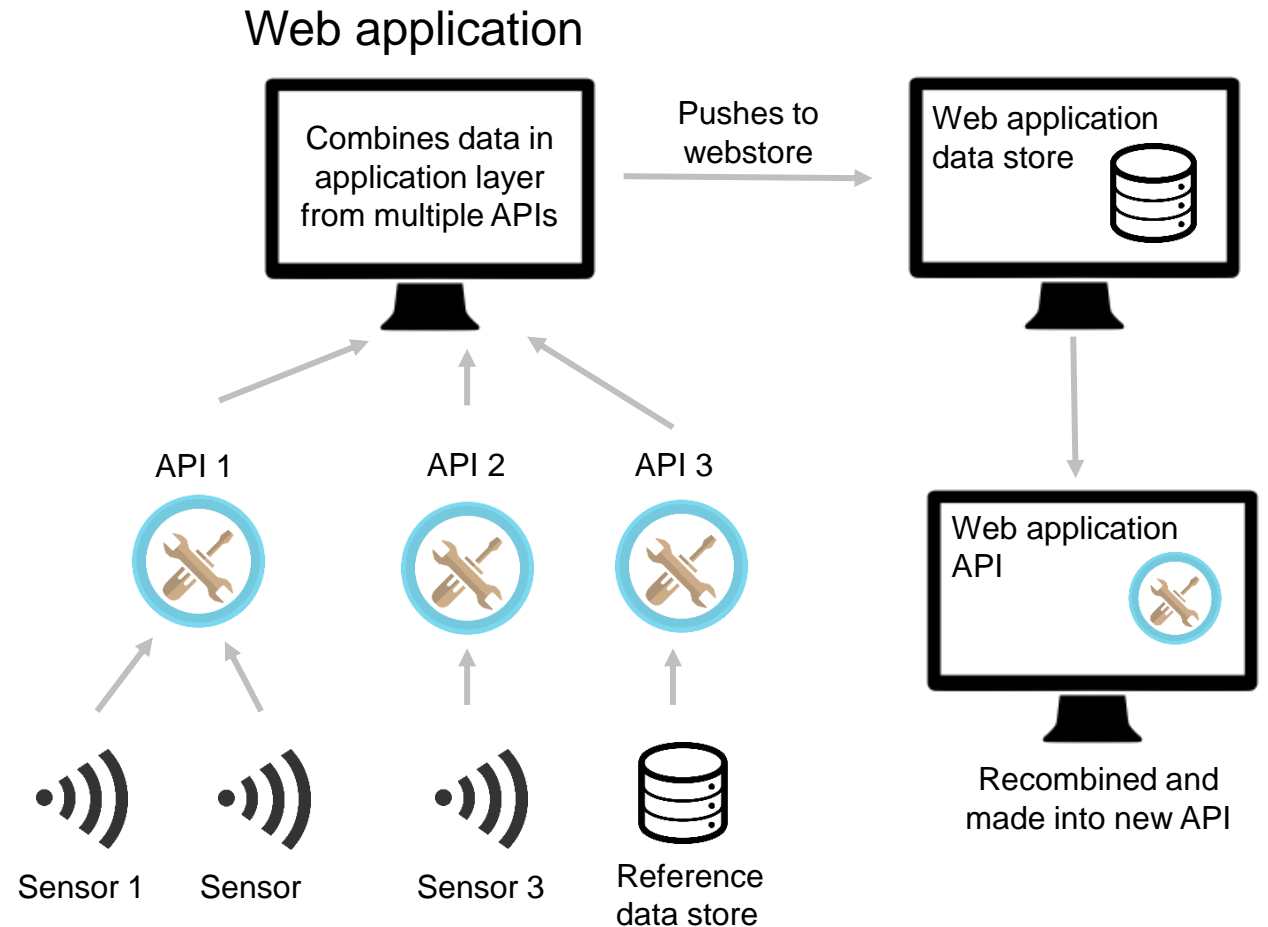Min and Max Profile | Data

# DI Use Case: Consume data from or check against a source
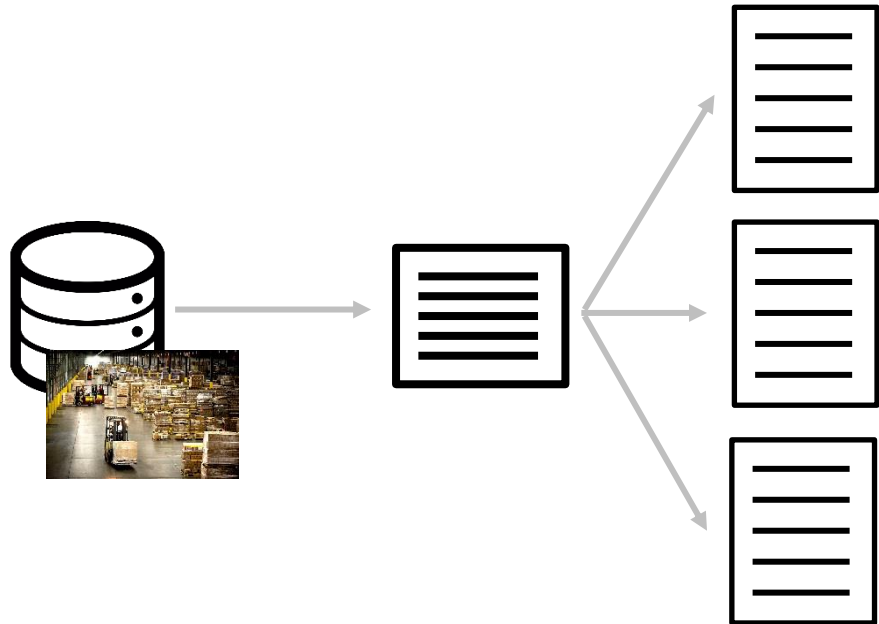
## E1 Simple yes/no check

Web application

Data store



Checks source

Yes! → Does something

Question?

No! → Does something else

## E2 Complex

Web application

Combines data in application layer from multiple APIs

Pushes to webstore →

Web application data store

API 1          API 2          API 3

Sensor 1   Sensor   Sensor 3   Reference data store

Web application API

Recombined and made into new API

# DI Use Case: Visualize data and KPIs, create dashboards

## 01 Centrally structured, codified, controlled, slow



Traditional data warehouse

Structured, controlled queries, views

Fixed, standard reports

## 02 Self-service, dynamic, ad hoc, fast



1 external source

CSV on desktop

Anything and everything

Locally trusted stores

# **Within** each DI use case, there are architectural choices, at a minimum…

Warehoused vs virtualization

Streaming versus batch

Kitchen sink vs lightly coupled commodity components

Central versus distributed

Low versus high volume

# **Between** DI use cases, earlier choices can restrict downstream choices

Warehoused vs virtualization

Streaming versus batch

Kitchen sink vs lightly coupled commodity components

Central versus distributed

Low versus high volume

Right now. We are making choices that will affect our downstream capacity and flexibility.

# Why we need a conscious data infrastructure strategy

- Improvements in data consistency and quality
- Faster, easier access to data
- Better controls and security
- Data sharing and interoperability between datasets
- Integrated data across departments
- Faster development of digital and web services
- Data analytics and more advanced data science
- New and novel data services

# Benefits to different groups

| Audience | Benefit |
|---|---|
| Developers and IT staff | Decrease in technical and development time to create applications, integrations and services |
| COIT and department budget staff | Decrease in costs for applications and services |
| Analysts, data users, ShareSF committee | More time and resources for conducting analysis and evaluation → better services and outcomes |
| Voters, program staff, executives | Better decisions and services |

# A possible reference architecture



| Governance, Policy, and Privacy Framework | Role Based Access Control | | | | | |
|---|---|---|---|---|---|---|
| | | Applications | Open Data Catalog | Apps and Visualizations | Business Intelligence Platforms | Etc. |
| | | Services | Managed APIs and API Gateway | Messaging Services | Metadata Platform & Services | Geo Services, etc |
| | | Sources for Distribution | Public Data Store with Public APIs | Streaming Data Service | Other Specialized Data Stores with APIs, etc. | |
| | | Dataset Development | Data Preparation, Cleansing and Review | | | |
| | | Connections and Transport | Connectors/Transport Layer | | | |
| | | Original Federated Sources | Oracle DB | SQL Server | Postgres | Etc. |

# Data, for the love of the City

THANK YOU
@datasf | datasf.org | datasf.org/blog

Data profiling

# APPENDIX

# DataSF aside: We've profiled every published dataset

# DataSF aside: And every published field

# DataSF aside: Profiling scripts are open source and building a dashboard so publishers can easily track