

Real-time Meatspace Data Science

Data Intelligence Conference

Jason Walsh
Data Engineer
Jason.Walsh@uphs.upenn.edu
@rightlag

Michael Becker
Sr. Data Scientist
Michael.Becker2@uphs.upenn.edu
@beckerfuffle

6/24/2017



Penn Medicine

Data Scientist or Data Janitor?



A blog about building better products

Stop Treating Your Data Scientist Like a Janitor

By [Archana Madhavan](#) | July 19 | [0 Comments](#)

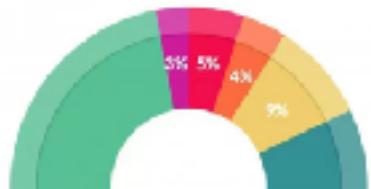


So you just hired someone to fill an opening for [the sexiest job of 2016](#), a data scientist at your up-and-coming Silicon Valley startup. What if, a few months down the line, you realize your data scientists are actually spending the majority of their time working on something decidedly... unsexy?

Stop Treating Your #DataScientist Like a Janitor

[CLICK TO TWEET](#)

Data cleanups aren't fun



What data scientists spend the most time doing

- Building training sets: 30%
- Cleaning and organizing data: 60%
- Collecting data sets: 10%

Data Scientist or Data Janitor?



A blog about building better products

Stop Treating Your Data Scientist Like a Janitor

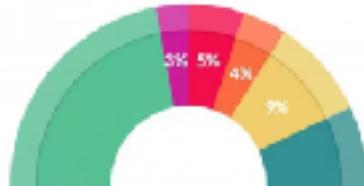
By Archana Madhavan | July 19 | 0 Comments



So you just hired someone to fill an opening for [the sexiest job of 2011](#) coming Silicon Valley startup. What if, a few months down the line, you're actually spending the majority of their time working on something di-

Stop Treating Your #DataScientist Like

Data cleanups aren't fun



What data scientists do:

- Building training sets
- Cleaning and organizing
- Collecting data sets
- Massaging data
- Modeling data

MAR 23, 2016 @ 09:33 AM 19,156 ▾

The Little Black Book of Billionaire Secrets

Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says



Gil Press, CONTRIBUTOR

I write about technology, entrepreneurs and innovation. [FULL BIO ▾](#)

Opinions expressed by Forbes Contributors are their own.

TWEET THIS

- Twitter icon: data scientists found that they spend most of their time massaging rather than mining or modeling data.
- Twitter icon: 76% of data scientists view data preparation as the least enjoyable part of their work

A new survey of data scientists found that they spend most of their time massaging rather than mining or modeling data. Still, most are happy with having [the sexiest job of the 21st century](#). The survey of about 80 data



Data Scientist or Data Janitor?

DNS

HOME

SEARCH

The New York Times

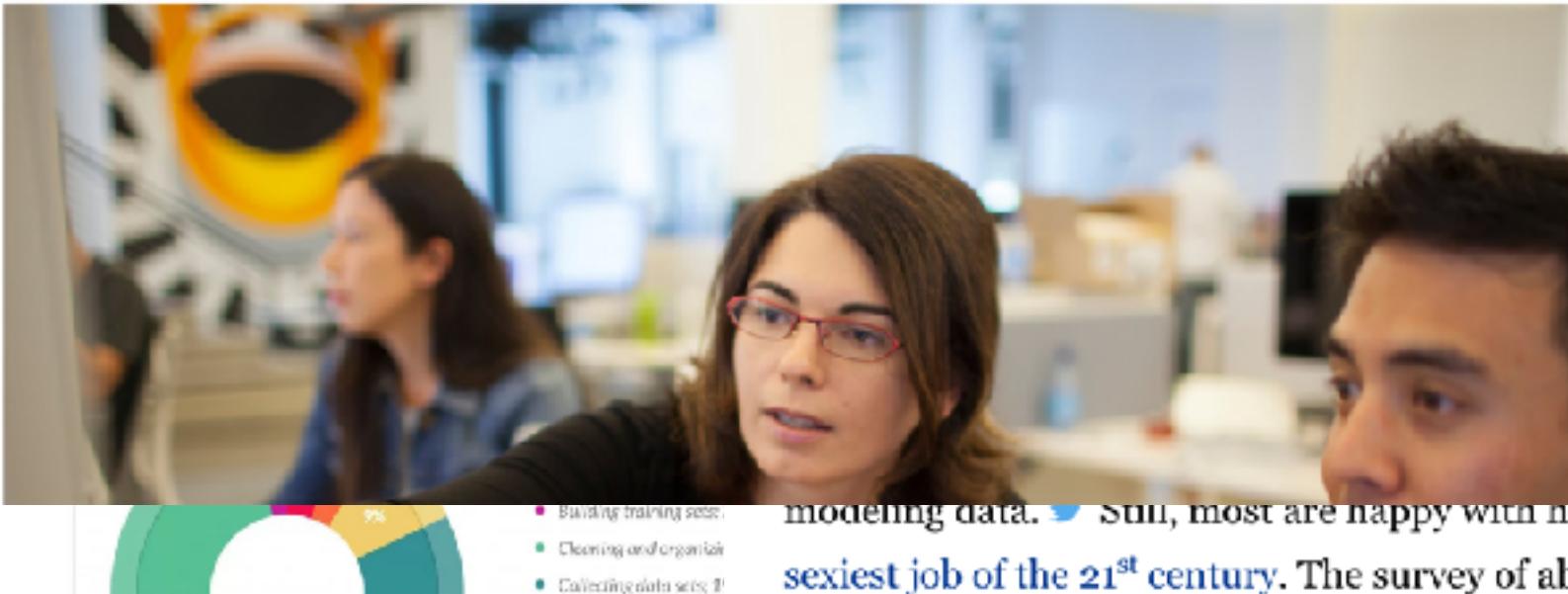
SUBSCRIBE NOW

LOG IN

TECHNOLOGY

For Big-Data Scientists, 'Janitor Work' Is Key Hurdle to Insights

By STEVE LOHR AUG. 17, 2014



- Building training sets
- Cleaning and organizing
- Collecting data sets

modeling data. Still, most are happy with having the sexiest job of the 21st century. The survey of about 80 data

Data Scientist or Data Janitor?



Big Data Borat

@BigDataBorat

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

Retweets

509

Likes

281



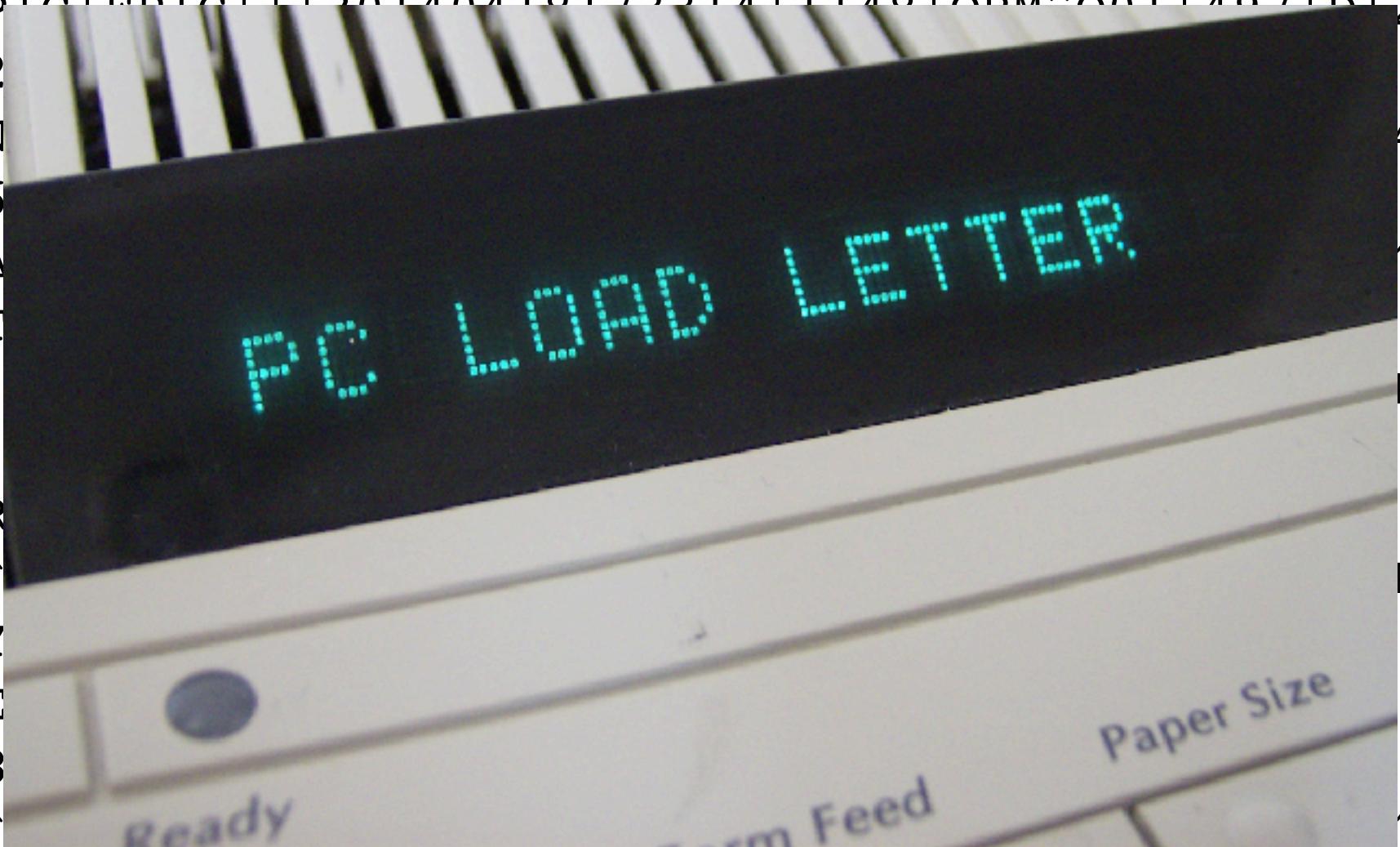
- Building training sets
- Cleaning and organizing
- Collecting data sets

modelling data. Still, most are happy with having the sexiest job of the 21st century. The survey of about 80 data

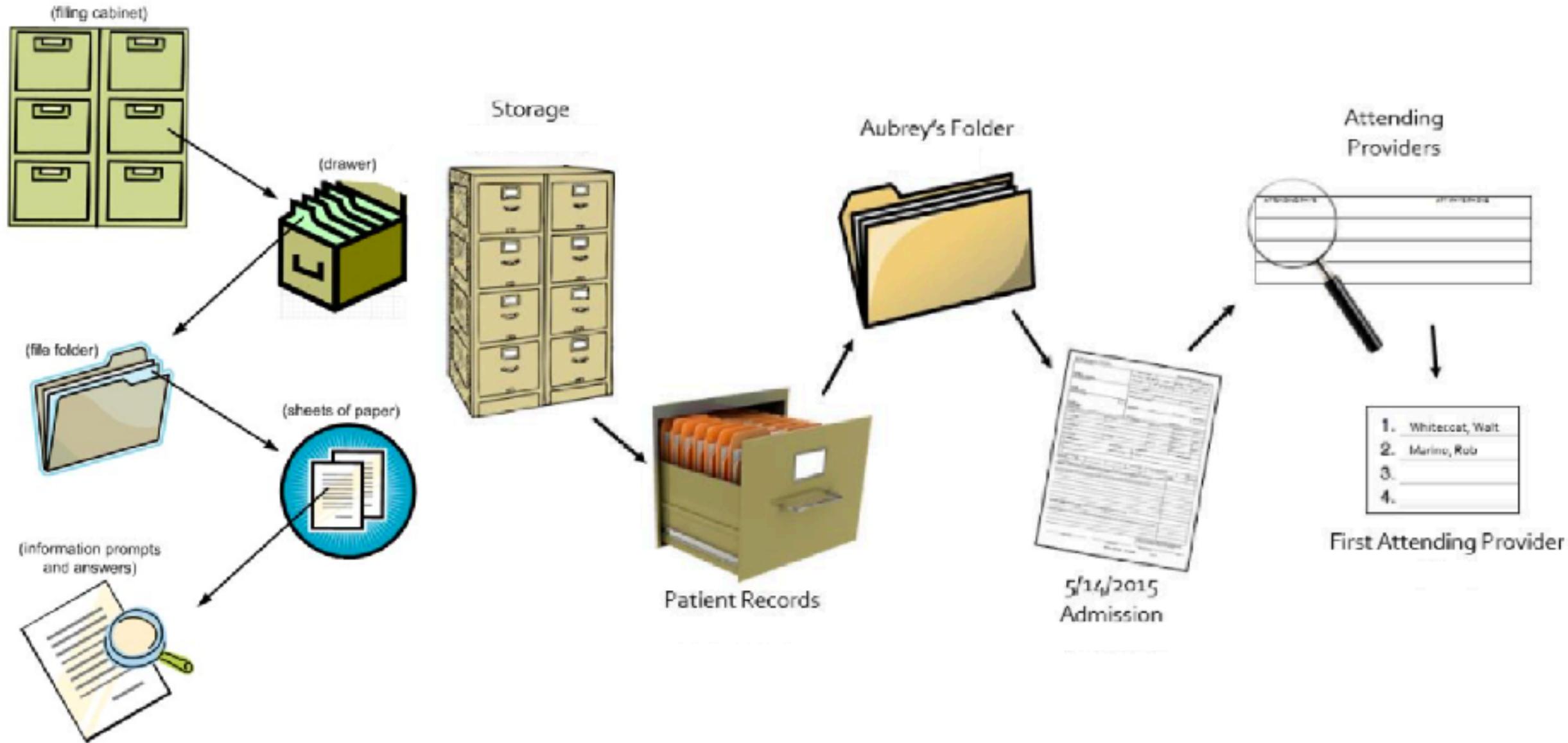
Health Care Data “Standards”

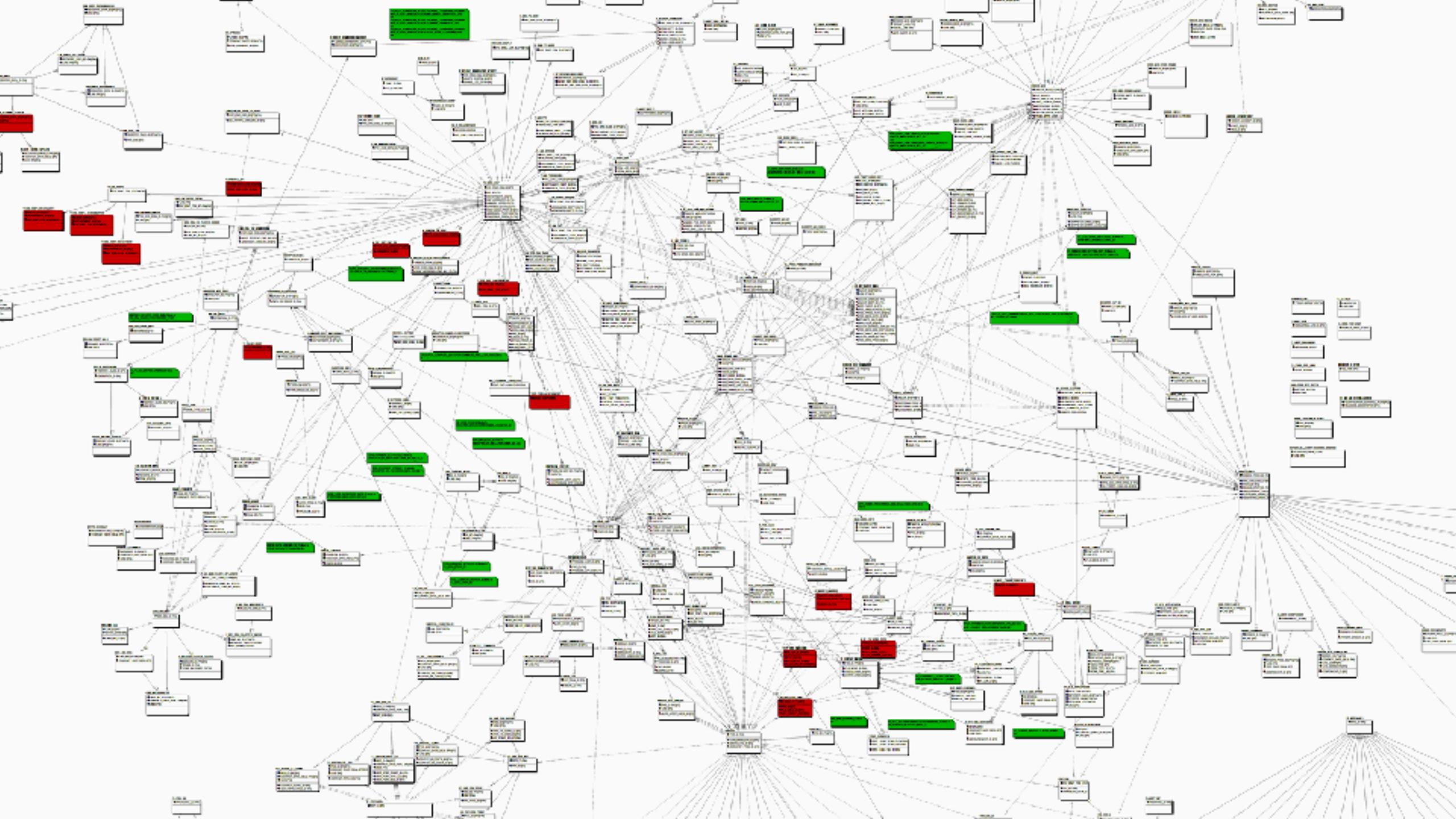
HL7

HL7

MSH|^~\&|EPIC|EPIC|||20140418172214|1148|OPM^001|407|P|2.3|| PID|
1||20891312|
S. HAMILTON|
(608)123-56|
FACILITY(EA|
HEALTH SYSTEM|
610613|||||
76543^EPC||
1148^PATTER|
1133^^222^
73610^X-RAY|
1173^MATTHE|
^^^20140418|
1148010^1A^|
FRACTURE|||
| AfrAm|505
4567|
|
^CARE|
| | | |
PIC|
|
PC|
|^I10|ANKLE

“Enterprise” Clinical Databases





My eyes! The goggles do nothing!



Health Care Data Issues

Inconsistent Data Types



Inconsistent Data Types

- ◆ unable to take temperature, patient just ate ice

Inconsistent Data Types

- ◆ unable to take temperature, patient just ate ice
- ◆ Pt c/o feeling hot, "like I have a fever". Deferred temp reassessment for 10min-pt recently sipped ice water

Inconsistent Data Types

- ◆ unable to take temperature, patient just ate ice
- ◆ Pt c/o feeling hot, "like I have a fever". Deferred temp reassessment for 10min-pt recently sipped ice water

An urge to eat non-nutritive substances, a condition called pica, **can** occur in **pregnant** women with low iron levels. On the other hand, **ice cravings**, known medically as pagophagia, don't necessarily **mean you** have a nutrient deficiency, which a blood test **can** diagnose. Apr 15, 2015

Inconsistent Data Types

- ◆ unable to take temperature, patient just ate ice
- ◆ Pt c/o feeling hot, "like I have a fever". Deferred temp reassessment for 10min-pt recently sipped ice water

An urge to eat non-nutritive substances, a condition called pica, **can** occur in **pregnant** women with low iron levels. On the other hand, **ice cravings**, known medically as pagophagia, don't necessarily **mean you** have a nutrient deficiency, which a blood test **can** diagnose. Apr 15, 2015

- ◆ unable to check temp with 3 different thermometers -- pt skin warm and dry -- bear hugger applied and will recheck temp soon
- ◆ BEAR HUGGER APPLIED

Inconsistent Data Types

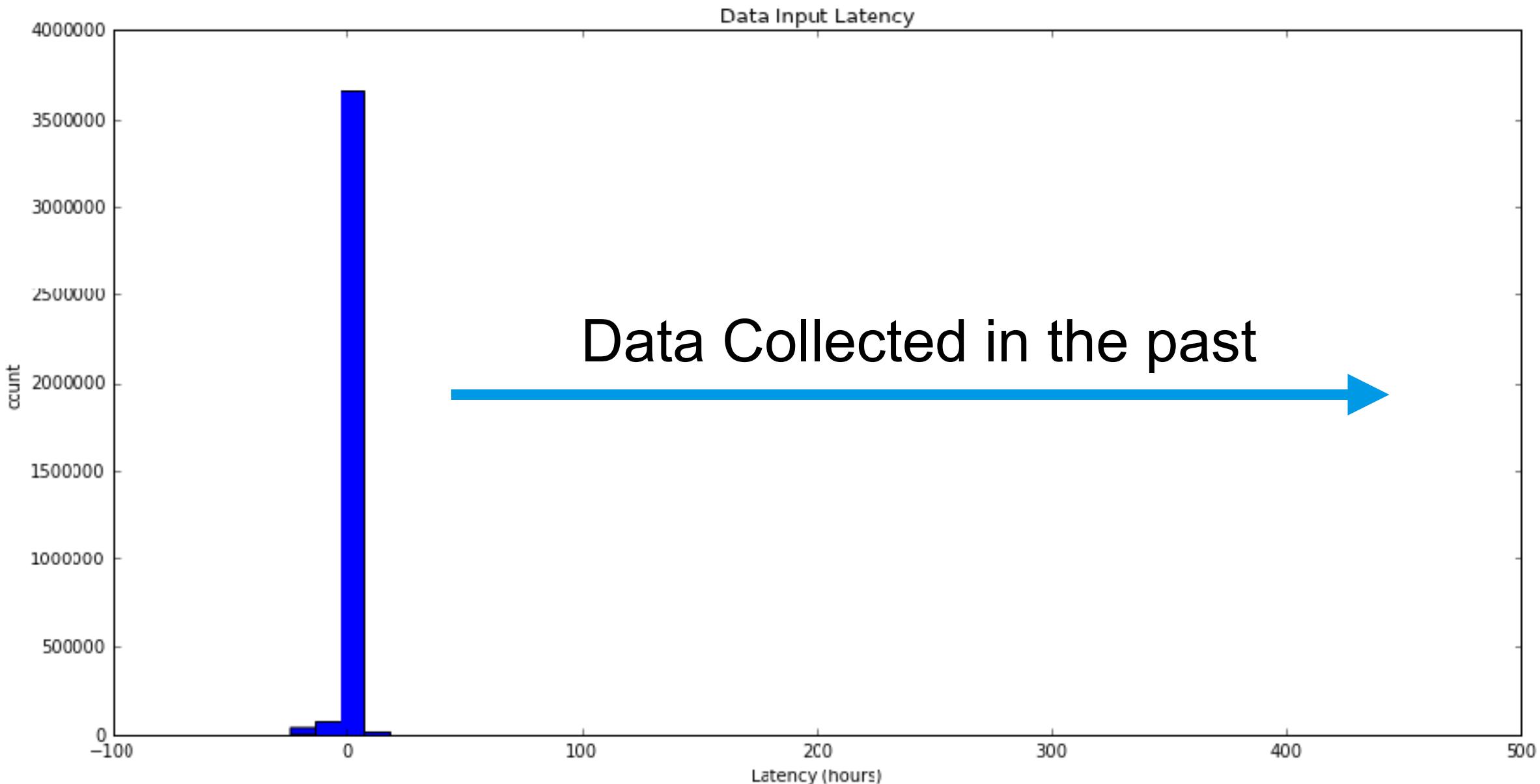
- ◆ unable to take temperature, patient just ate ice
- ◆ Pt c/o feeling hot, "like I have a fever". Deferred temp reassessment for 10min-pt recently sipped ice water

An urge to eat non-nutritive substances, a condition called pica, **can** occur in **pregnant** women with low iron levels. On the other hand, **ice cravings**, known medically as pagophagia, don't necessarily **mean you** have a nutrient deficiency, which a blood test **can** diagnose. Apr 15, 2015

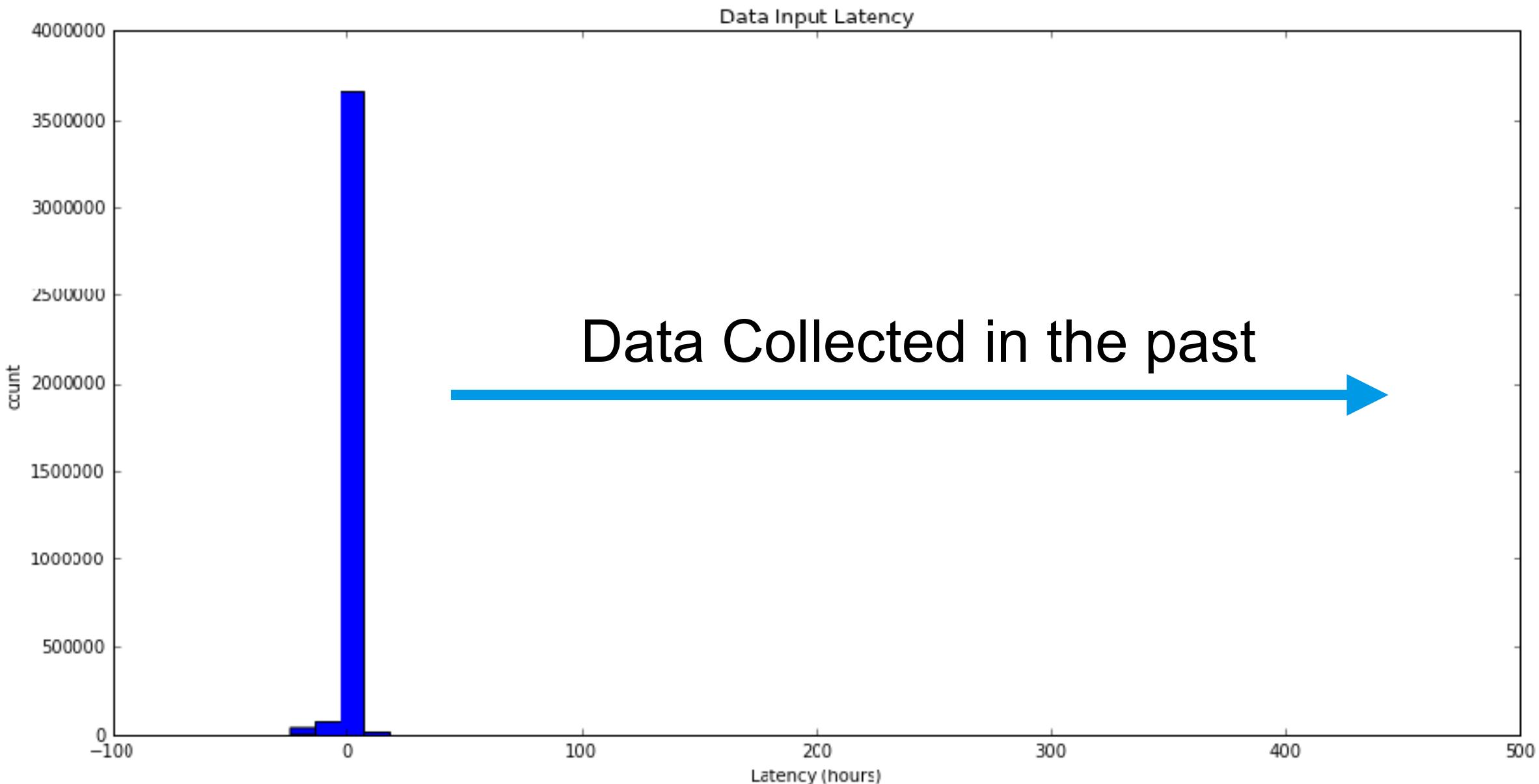
- ◆ unable to check temp with 3 different thermometers -- pt skin warm and dry -- bear hugger applied and will recheck temp soon
- ◆ BEAR HUGGER APPLIED



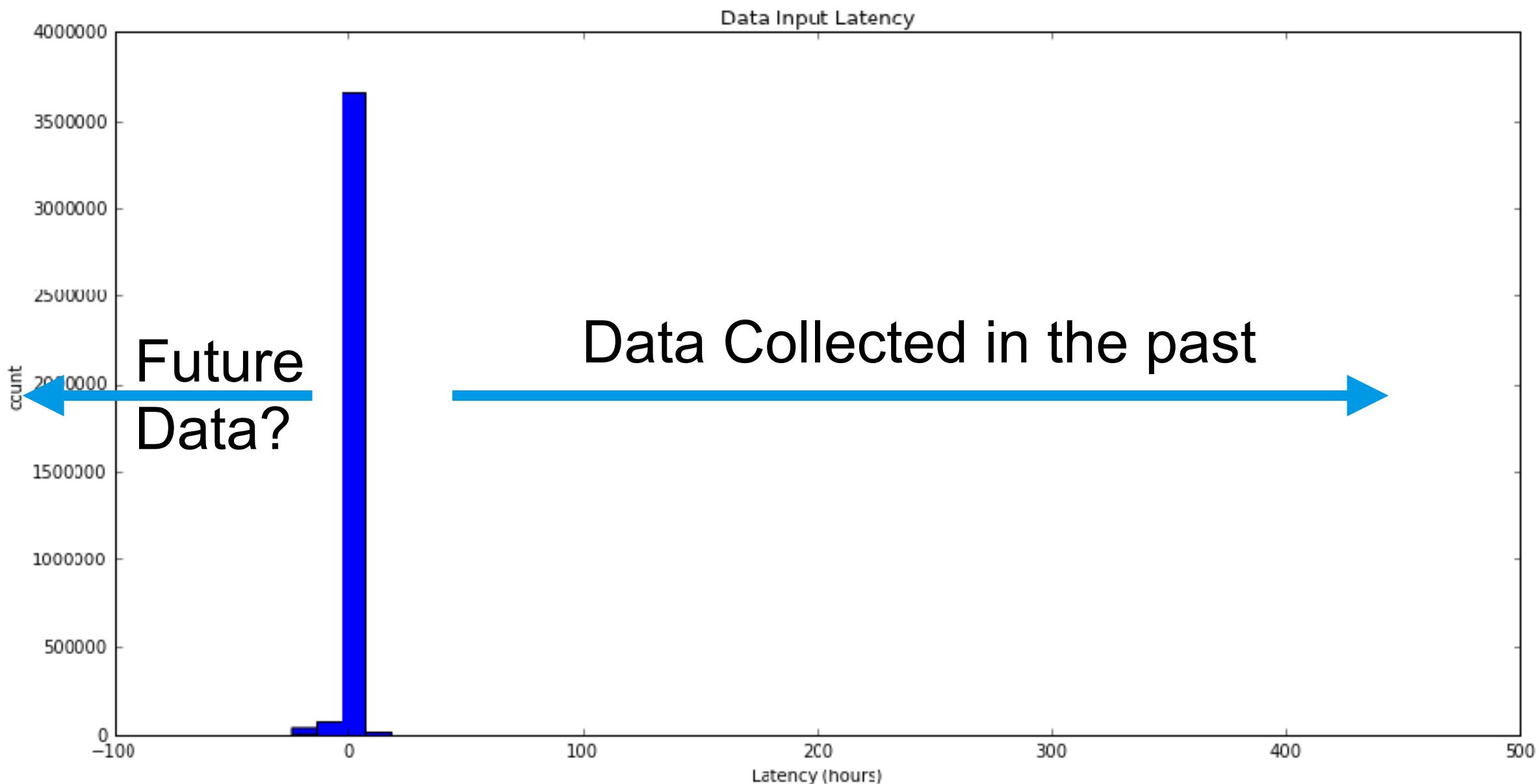
Future Data



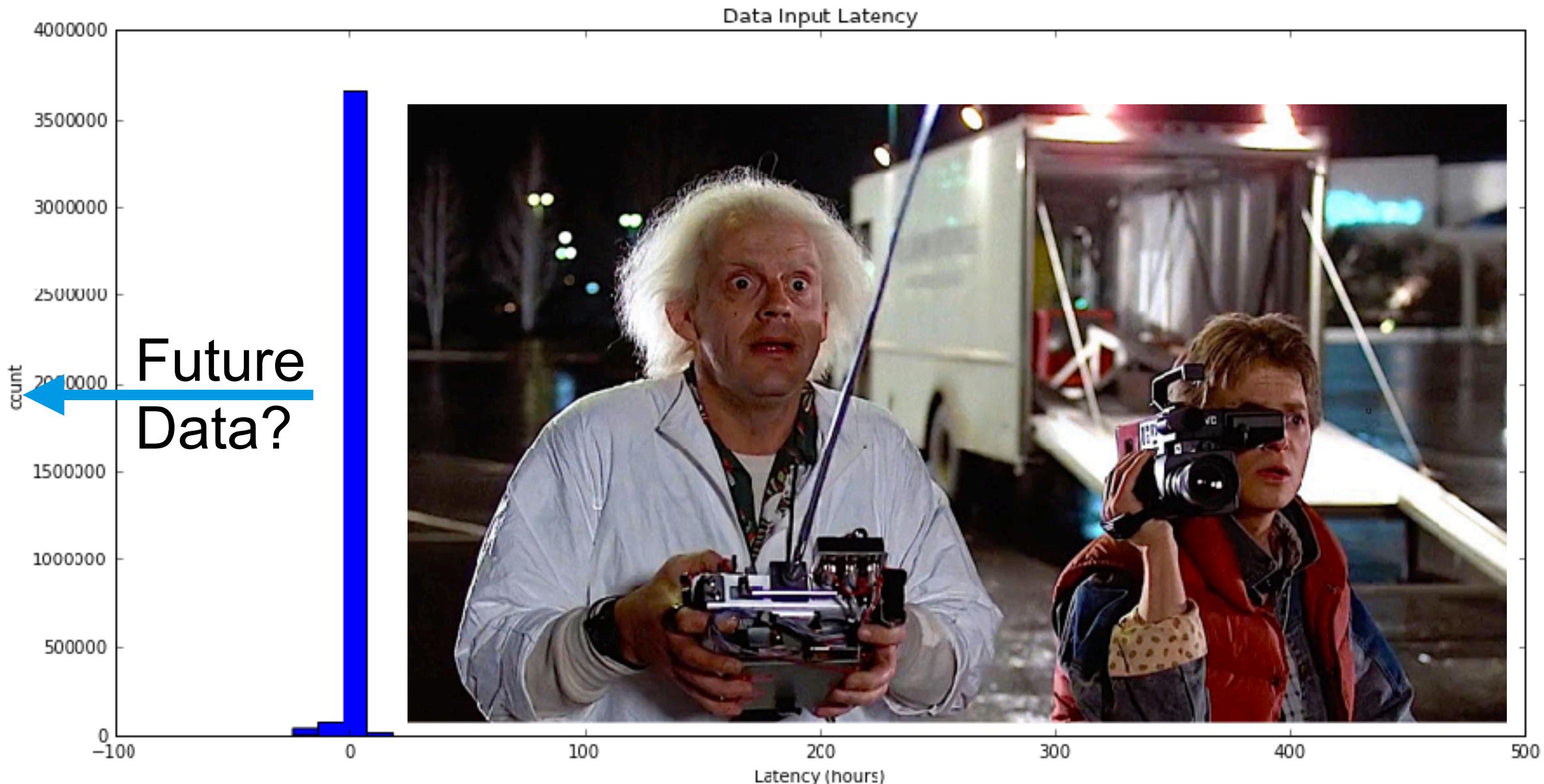
Future Data



Future Data



Future Data



More Data Issues

- ◆ Multiple identifiers for the same patient

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data
- ◆ Data collection biases

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data
- ◆ Data collection biases
- ◆ Inconsistent database models

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data
- ◆ Data collection biases
- ◆ Inconsistent database models
- ◆ Mostly human entered data... in a time crunch data is lost

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data
- ◆ Data collection biases
- ◆ Inconsistent database models
- ◆ Mostly human entered data... in a time crunch data is lost
- ◆ Little validation of input

More Data Issues

- ◆ Multiple identifiers for the same patient
- ◆ Throwing out historic data
- ◆ Data collection biases
- ◆ Inconsistent database models
- ◆ Mostly human entered data... in a time crunch data is lost
- ◆ Little validation of input

Little understanding of business impact

Penn Signals

Penn Signals

- ◆ A single data schema for data in stream and at rest

Penn Signals

- ◆ A single data schema for data in stream and at rest
- ◆ Enables fast transition from modeling to applications

Penn Signals

- ◆ A single data schema for data in stream and at rest
 - ◆ Enables fast transition from modeling to applications
- ◆ Storing event stream allows you to recreate the state of a patient at any time in their stay

Penn Signals

- ◆ A single data schema for data in stream and at rest
 - ◆ Enables fast transition from modeling to applications
- ◆ Storing event stream allows you to recreate the state of a patient at any time in their stay
 - ◆ Allows for deterministic modeling

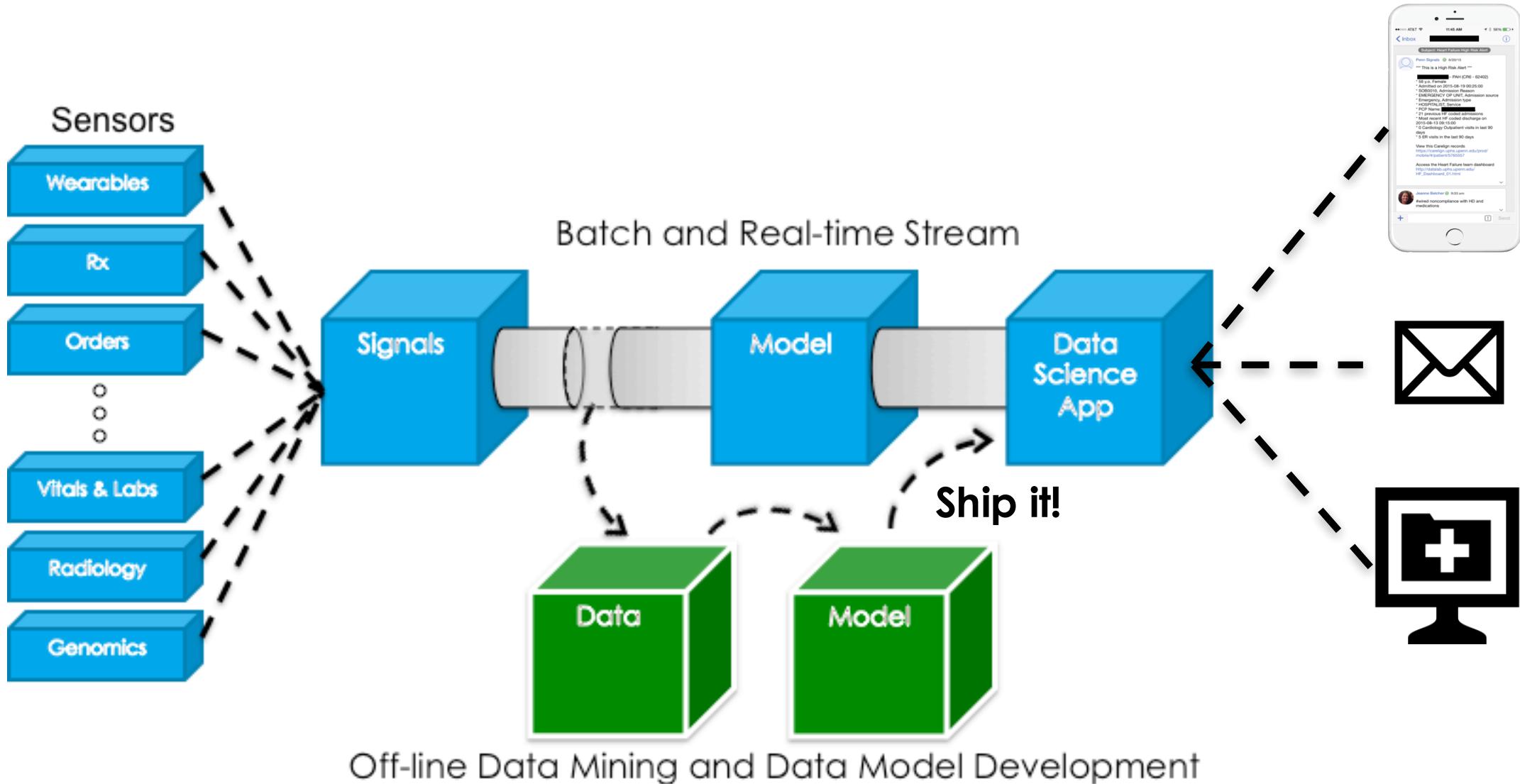
Penn Signals

- ◆ A single data schema for data in stream and at rest
 - ◆ Enables fast transition from modeling to applications
- ◆ Storing event stream allows you to recreate the state of a patient at any time in their stay
 - ◆ Allows for deterministic modeling
 - ◆ Crucial for real-time applications

Penn Signals

- ◆ A single data schema for data in stream and at rest
 - ◆ Enables fast transition from modeling to applications
- ◆ Storing event stream allows you to recreate the state of a patient at any time in their stay
 - ◆ Allows for deterministic modeling
 - ◆ Crucial for real-time applications
- ◆ Type enforcement

Penn Signals Architecture



Penn Signals Stack



Nomad



ØMQ

Spark

The Consul logo features a magenta stylized 'C' with a cluster of dots inside it, followed by the word "Consul" in black.

Vault

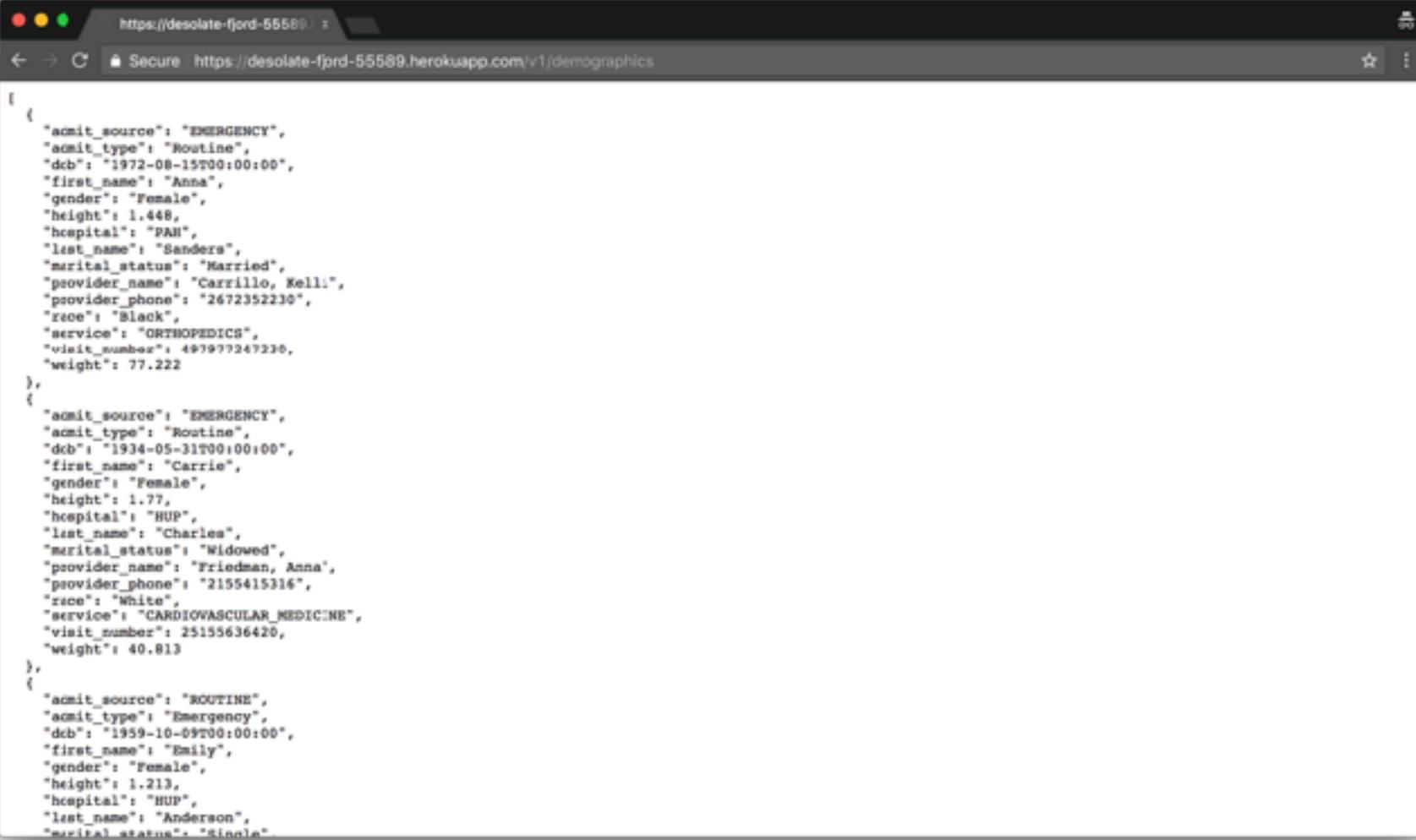
The Jupyter logo features a stylized orange and grey circular icon followed by the word "jupyter" in a sans-serif font.

Upstream Resources

- ◆ Enterprise data warehouses
- ◆ Web services
- ◆ Streaming telemetry data

fuzzy-engine

- ◆ <https://github.com/pennsignals/fuzzy-engine>
- ◆ returns JSON formatted HTTP responses containing randomized patient data



A screenshot of a web browser window displaying a JSON response. The URL in the address bar is `https://desolate-fjord-55589.herokuapp.com/v1/demographics`. The JSON data consists of three objects, each representing a patient record:

```
[  
  {  
    "admit_source": "EMERGENCY",  
    "admit_type": "Routine",  
    "dob": "1972-08-15T00:00:00",  
    "first_name": "Anna",  
    "gender": "Female",  
    "height": 1.448,  
    "hcopital": "PAH",  
    "last_name": "Sanders",  
    "marital_status": "Married",  
    "provider_name": "Carrillo, Kelli",  
    "provider_phone": "2672352230",  
    "race": "Black",  
    "service": "ORTHOPEDICS",  
    "visit_number": 497977247230,  
    "weight": 77.222  
  },  
  {  
    "admit_source": "EMERGENCY",  
    "admit_type": "Routine",  
    "dob": "1934-05-31T00:00:00",  
    "first_name": "Carrie",  
    "gender": "Female",  
    "height": 1.77,  
    "hcopital": "HUP",  
    "last_name": "Charles",  
    "marital_status": "Widowed",  
    "provider_name": "Friedman, Anna",  
    "provider_phone": "2155415316",  
    "race": "White",  
    "service": "CARDIOVASCULAR_MEDICINE",  
    "visit_number": 25155636420,  
    "weight": 40.813  
  },  
  {  
    "admit_source": "ROUTINE",  
    "admit_type": "Emergency",  
    "dob": "1959-10-09T00:00:00",  
    "first_name": "Emily",  
    "gender": "Female",  
    "height": 1.213,  
    "hcopital": "HUP",  
    "last_name": "Anderson",  
    "marital_status": "Single"  
  }]
```

Can we validate upstream data?

```
{  
  "id": 1,  
  "name": "A green door",  
  "price": 12.50,  
  "tags": [ "home", "green" ]  
}
```



JSON Schema Vocabulary

```
{  
  "$schema": "http://json-schema.org/draft-04/schema#",  
  "title": "Product",  
  "description": "A product from Acme's catalog",  
  "type": "object",  
  "properties": {  
    "id": {  
      "description": "The unique identifier for a product",  
      "type": "integer"  
    },  
    "name": {  
      "description": "Name of the product",  
      "type": "string"  
    },  
    "price": {  
      "type": "number",  
      "minimum": 0,  
      "exclusiveMinimum": true  
    },  
    "tags": {  
      "type": "array",  
      "items": {  
        "type": "string"  
      },  
      "minItems": 1,  
      "uniqueItems": true  
    }  
  },  
  "required": ["id", "name", "price"]  
}
```

JSON Schema Validation Keywords

strings

- maxLength
- minLength
- pattern

arrays

- items
- additionalItems
- maxItems
- minItems
- uniqueItems

That's cool, but are there documents that describe web services?

YES!



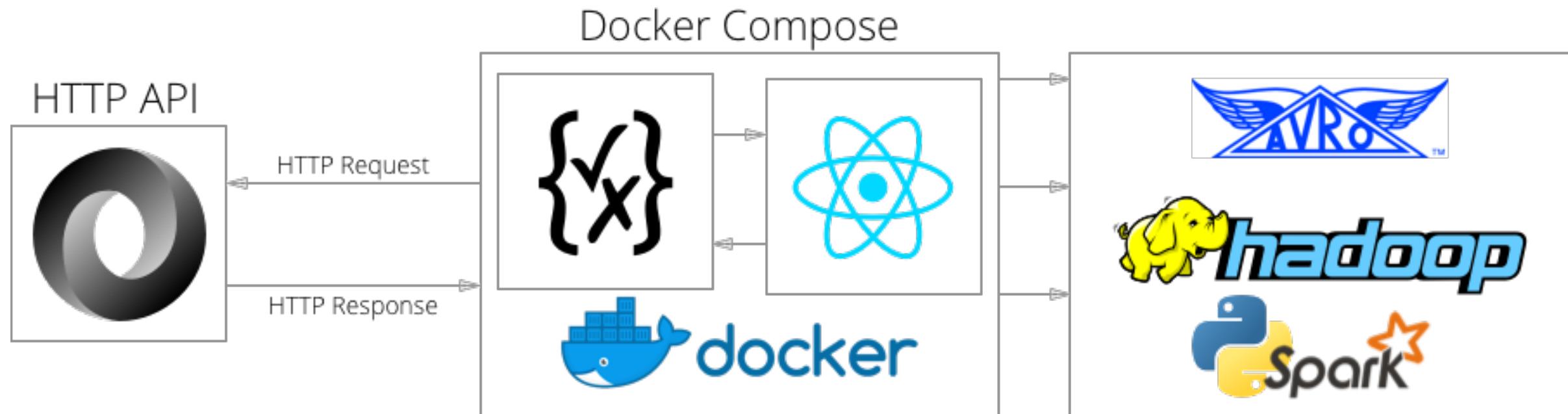


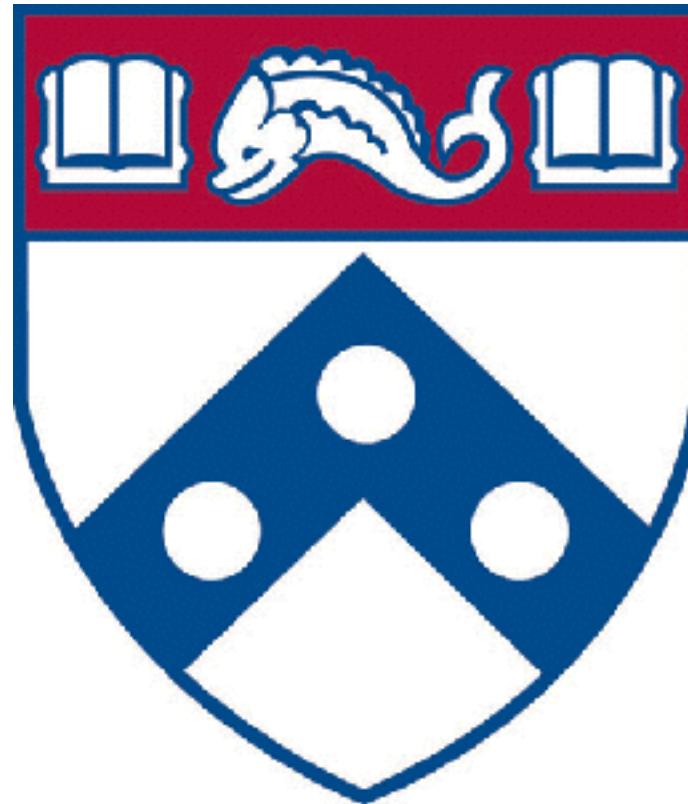
<http://swagger.io/>

Putting it all together



Process





Data Scientist or Data Janitor?



Big Data Borat

@BigDataBorat

In Data Science, 80% of time spent prepare data, 20% of time spent complain about need for prepare data.

Retweets

509

Likes

281



About us

- ◆ Twitter @beckerfuffle & @rightlag
- ◆ Follow our team @PennDataScience
- ◆ Updates on Penn Signals can be found at predictivehealthcare.pennmedicine.org