

# Data Workshop: Tackling RWD Integration Challenges

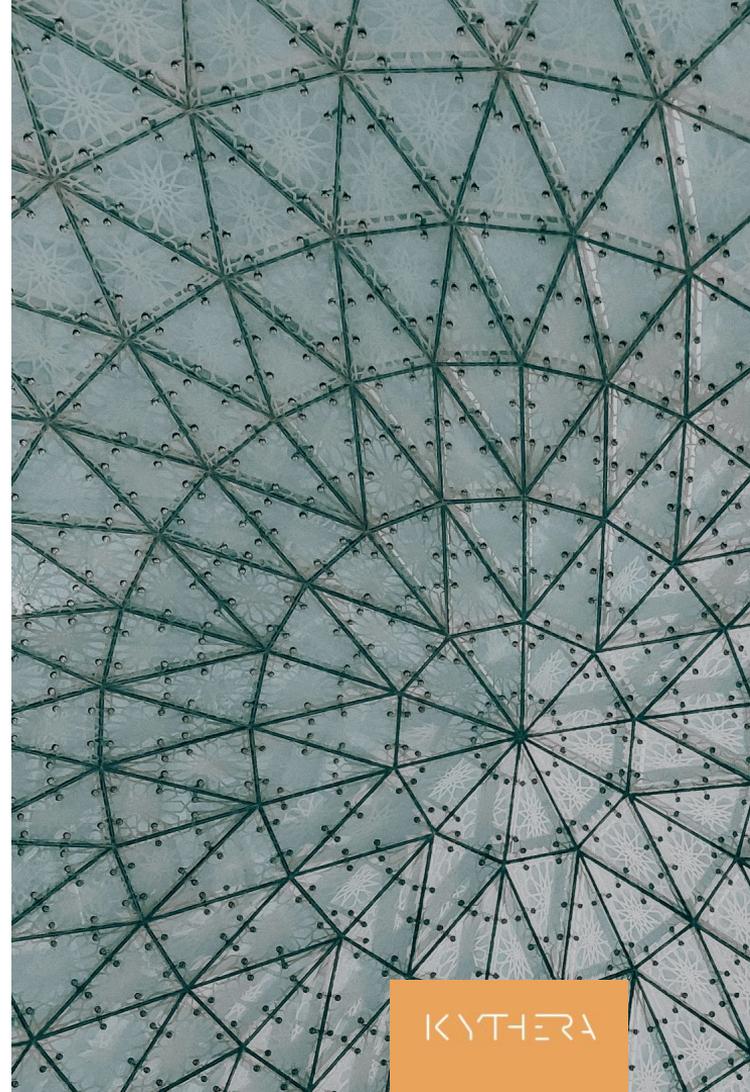
KYTERA  
DECIPHERING HEALTHCARE

 NOVARTIS

# Kythera Labs' Mission

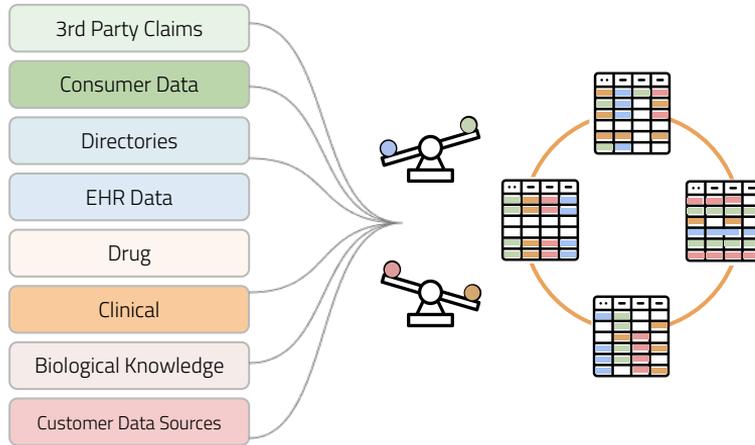
---

Reduce the uncertainty in the use  
of real-world data.

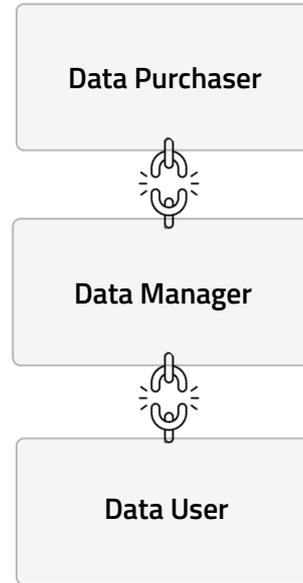


# Understanding the Real-World Data Integration Problem

1. Data source **bias, fragmentation, and overlap**



2. Common **disconnect** between data purchasers, managers, and users

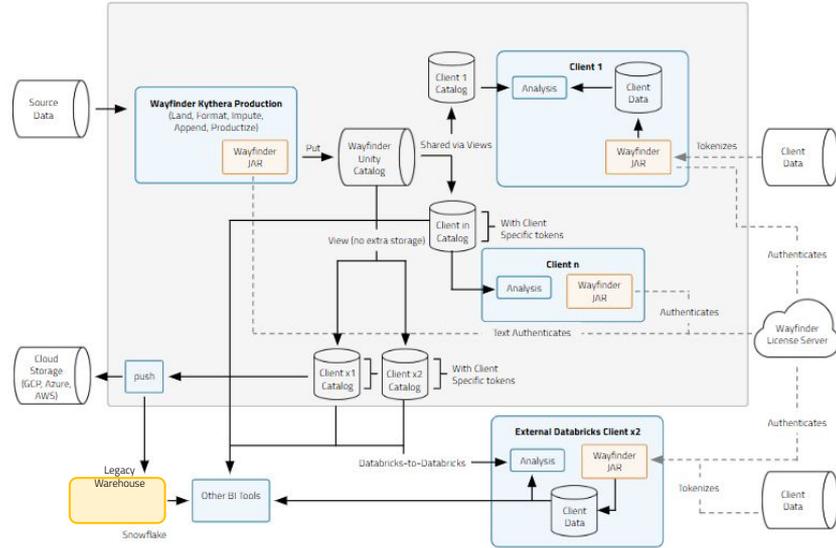


# Understanding the Real-World Data Integration Problem

## 3. Seeing only **parts** of the picture



## 4. A **simpler** solution is highly **complex**



# The True Costs of Real-World Data

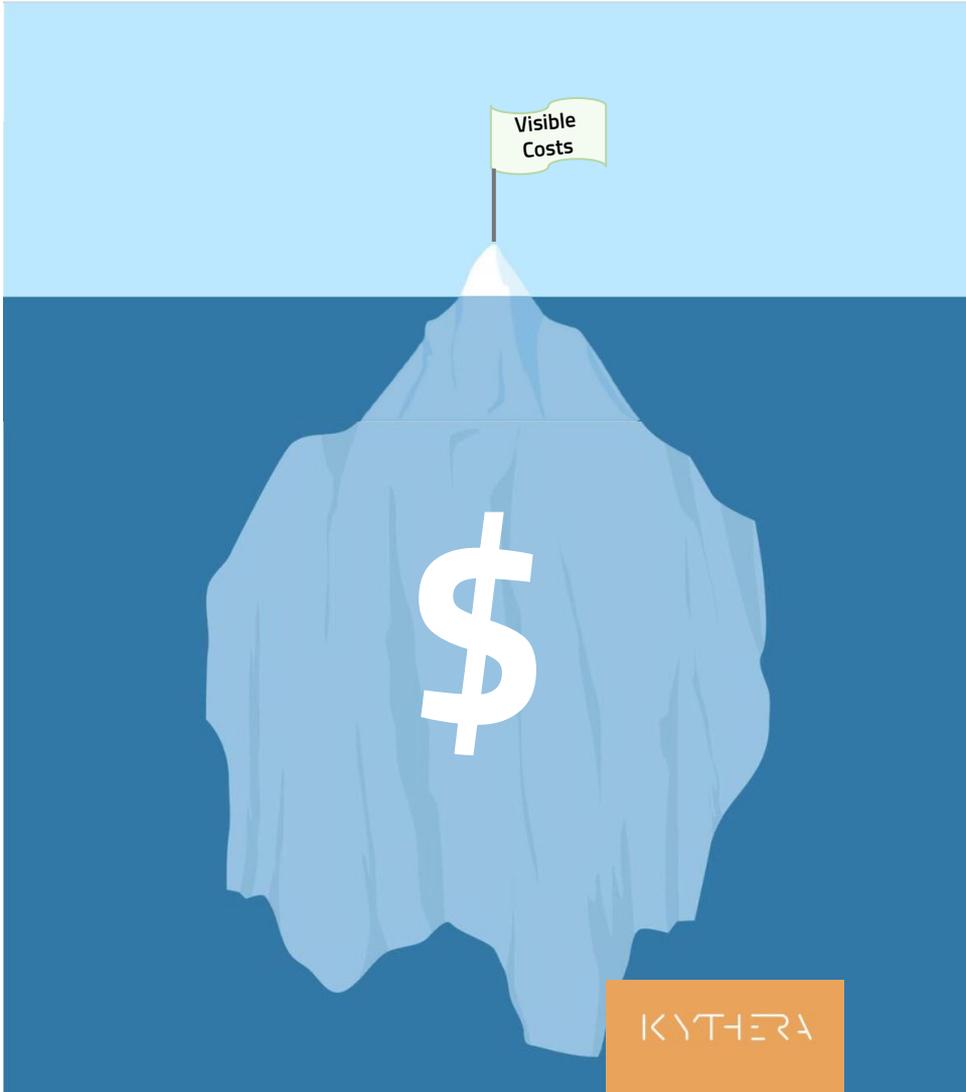
---

## Visible Costs

- Data Licensing Fees
- Integration Tools
- Internal Resources
- External Resources

## Not-So-Visible Costs

- Time delays
- Insufficient clean-up results
- Redundant purchases
- Opportunity costs
- Failed use cases
- Missed regulatory insights
- Misaligned commercialization



# Relevant Questions To Ask Yourself

---



How confident are you in the **completeness** of patient journeys in your most-used datasets?



Have you ever run an analysis, only to later find out the dataset missed a significant portion of the **target population**?



Do you know how many different **third-party datasets** your team licenses? How many are linked?



When was the last time a **tokenization issue** disrupted your insight delivery or caused delays?



What percent of your RWD budget do you think is spent **managing data**, not using it?



# High Level Outline

---

- Use case can benefit from combination of 3rd party data sets to improve accuracy and completeness of analyses
- 3rd party data often de-identified for privacy and compliance
- De-identification introduces new challenges, patient mastering addresses these
  - Patient duplication within a single dataset
    - Gaps in patient history
    - **Method: Identify “suspect tokens”**
    - **Method: Match tokens and dedupe**
  - Token scheme mismatches between datasets
    - Token “Drop-off” when joining
    - **Method: Token bridging**
- Assessing Completeness of data to support analyses
  - **Method: Effective overlap analyses**
  - **Method: Identify unique uplift of a data set**

# Closed & Open Claims

Patient Encounter Data:

## Closed Claims

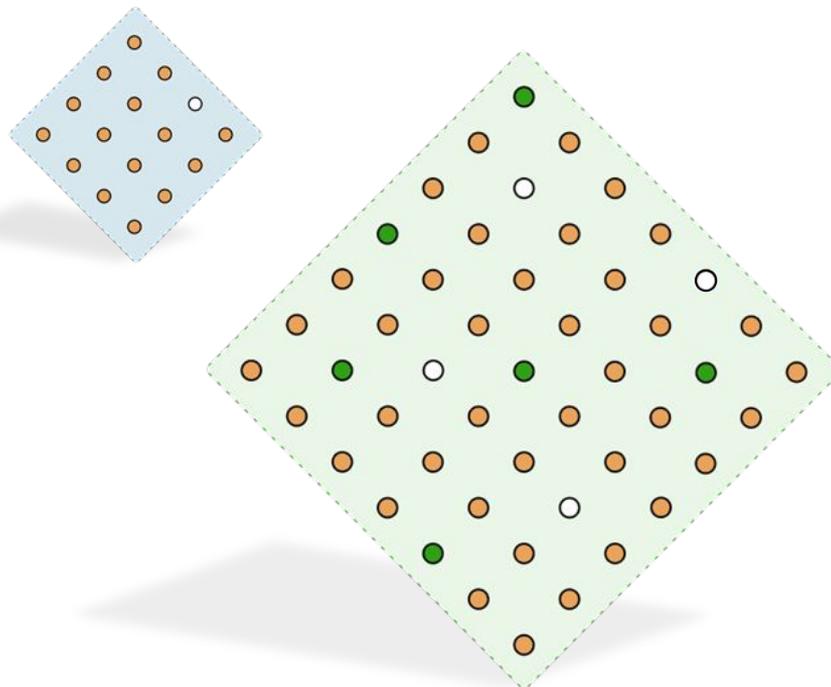
### PROS

- Adjudicated claims data
- Full records of care
- High rate of data completion

### CONS

- Expensive
- Payer-biased
- Region-biased
- Latency
- Short window of visibility
- Small sample size

- Encounter Data
- Kythera Imputed data
- Missing data



Patient Encounter Data:

## Open Claims

### PROS

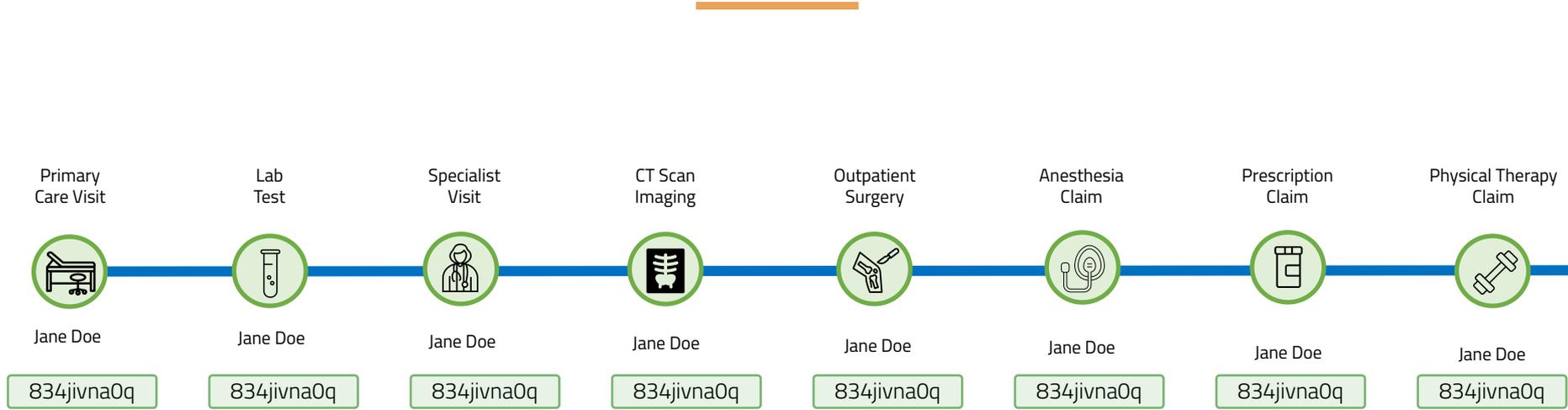
- Less costly
- Frequently updated
- Longer window of visibility
- Much larger sample size
- Nationally representative

### CONS

- ~~Provider biased~~
- ~~Input errors~~
- ~~Missingness~~
- ~~Patient duplication~~

# Patient Journey Ideal State

One Patient : One Journey



# Tokenization Intro

- De-Identification
  - John Q. Smith → “M2s\$\*znM6A3YYu7T”
  - Increase privacy
  - Increase acceptable use cases for data
- Security
  - Prevent leaks of PII
  - Enable collaboration between organizations that require patient identifiers



# Tokenization Can Introduce New Errors

Challenges with tokenization for the patient journey

Accurate Token Match

**One Person : One Token**



s6t8bdtuk8s6fgh1s6r8t

False Negative Token Match

**One Person : Multiple Tokens**

v0ertgo6vqapirfpa89a



834jivna0q9quinfv08sd

False Positive Token Match

**Multiple Persons : One Token**

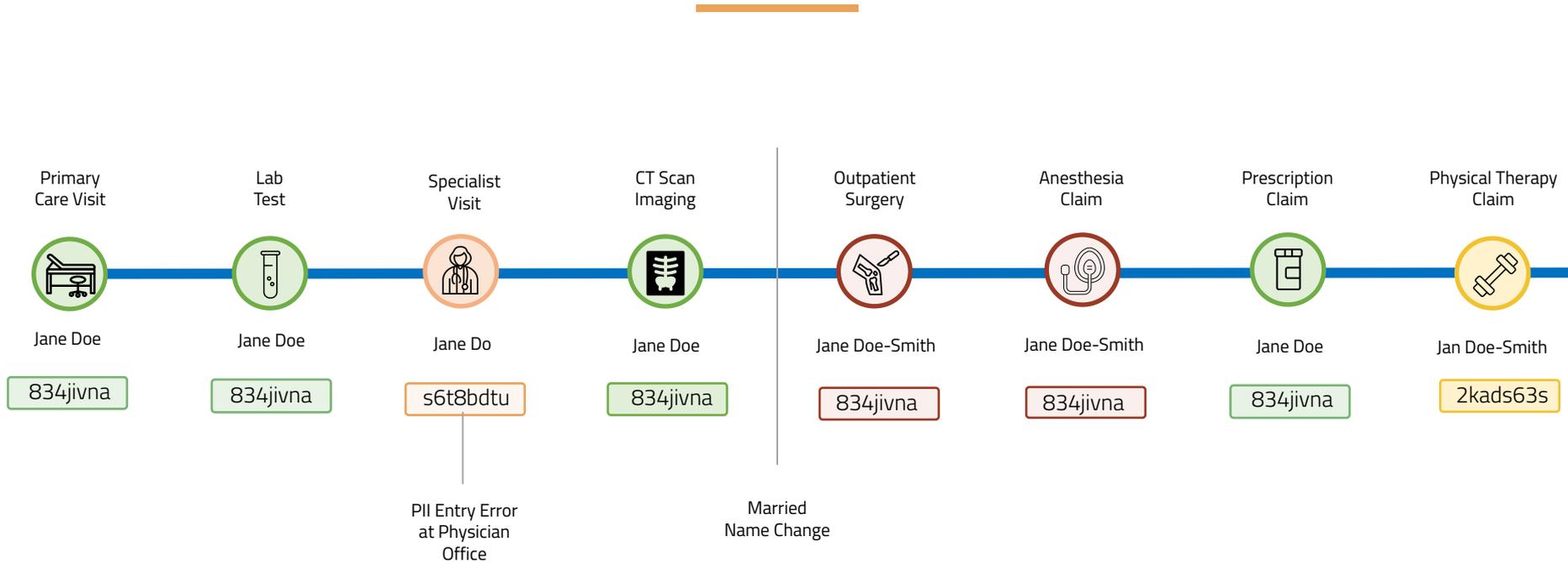


s6t8bdtuk8s6fgh1s6r8t



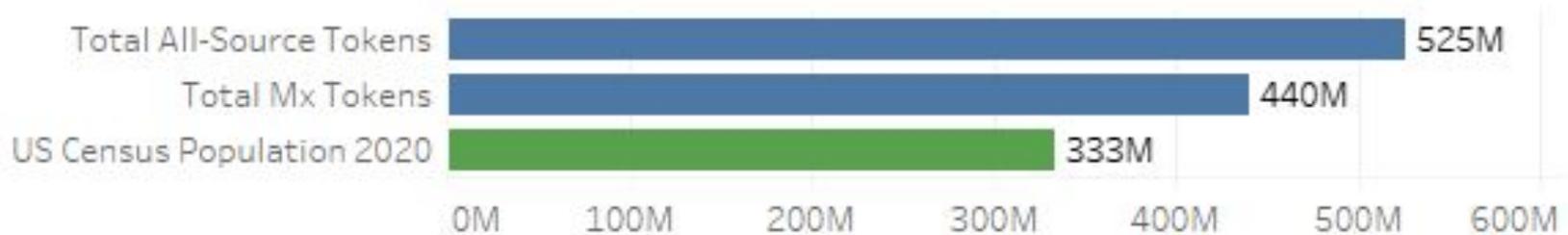
# Patient Journey Fragmented State

One Patient : Multiple Tokens



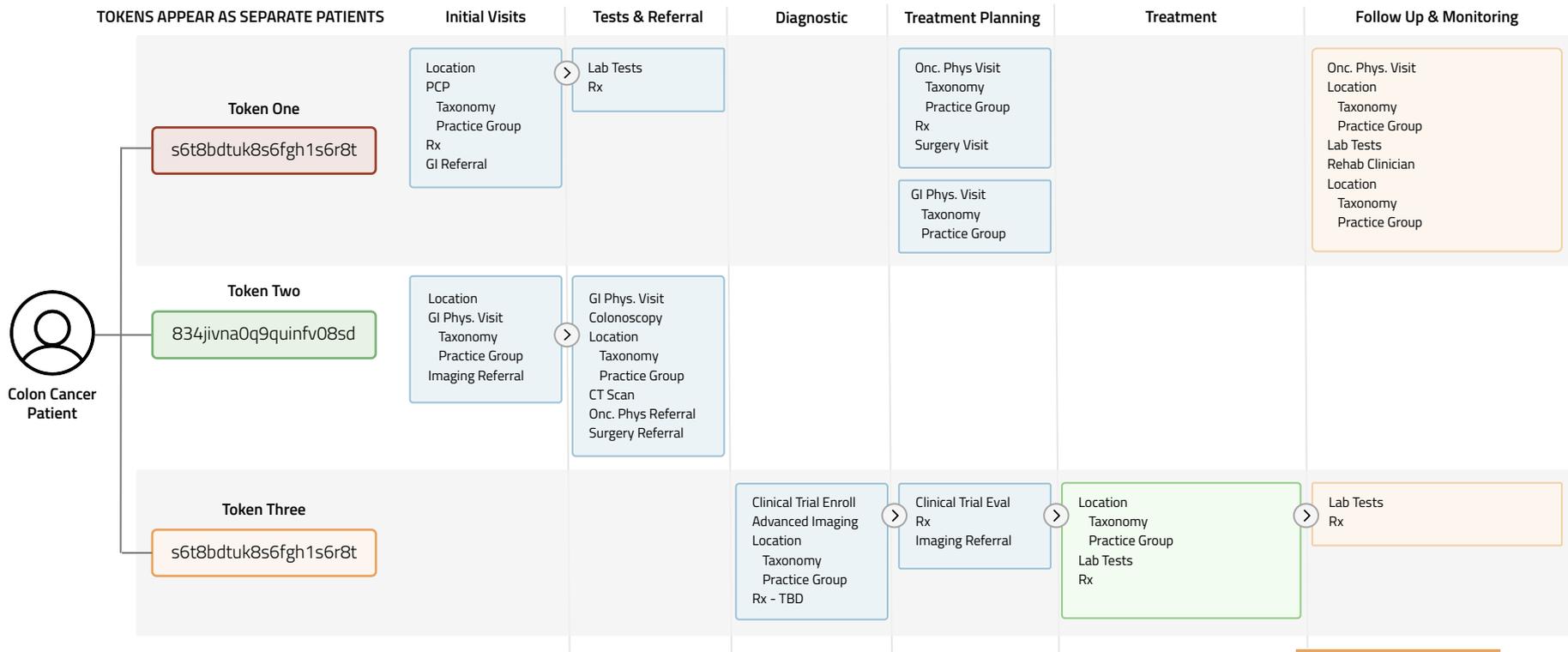
# Extent of the Problem: Single Data Source

- From a single vendor of clearinghouse claims (Mx and Rx) Kythera sees over **520M** distinct token combinations (token1 and token2 together), which is clearly impossible for a dataset containing only information on US individuals if those are supposed to be distinct patients
- Considering the medical submits only there are **440M** distinct token combinations – better, but still far too high



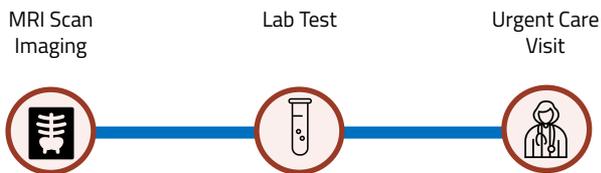
# Token Misalignment

Multiple tokens result in disparate, unreliable, inaccurate event capture



# Suspect Token Identification

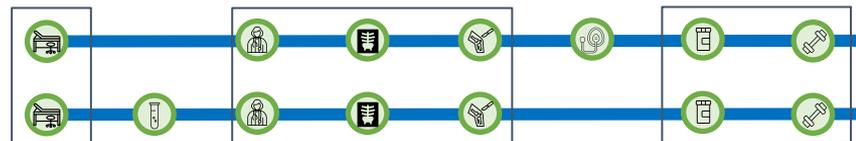
- Suspect tokens are tokens that are statistically likely to represent only a fragment of a complete patient journey.
- Suspect tokens may appear for many reasons, such as by a patient changing their name or an input error at the provider level.



Age:	63 years of age
Gender:	Females
Diagnosis Code:	M4726 (Degenerative Spine Condition)
Procedure Code:	72149 (MRI)
Number of patients in the Kythera claims set who present with both this primary diagnosis and procedure	1,119 patients
Mean chronological count	166 (for the 1119 patients, the average number of claims generated by the patients before the MRI for spondylolysis of the lumbar region is 165).
Standard Deviation of the chronological count	220
Chronological count for the lowest 1% of the population	149

# Patient Dedupe Methodology #1

- Train a statistical ML model using closed claims
  - Likelihood of patients appearing at same provider, same day, same procedure
  - Prior distribution on likelihood of patients sharing same chronic diagnoses and comorbidities
- Identify candidate token matches using similarity measures
- Evaluate the quality of the candidate match

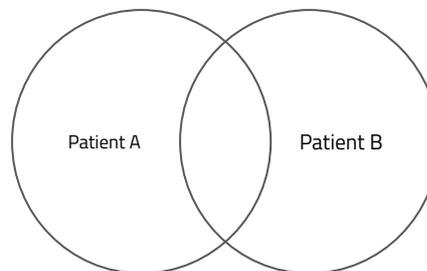


## Patient A:

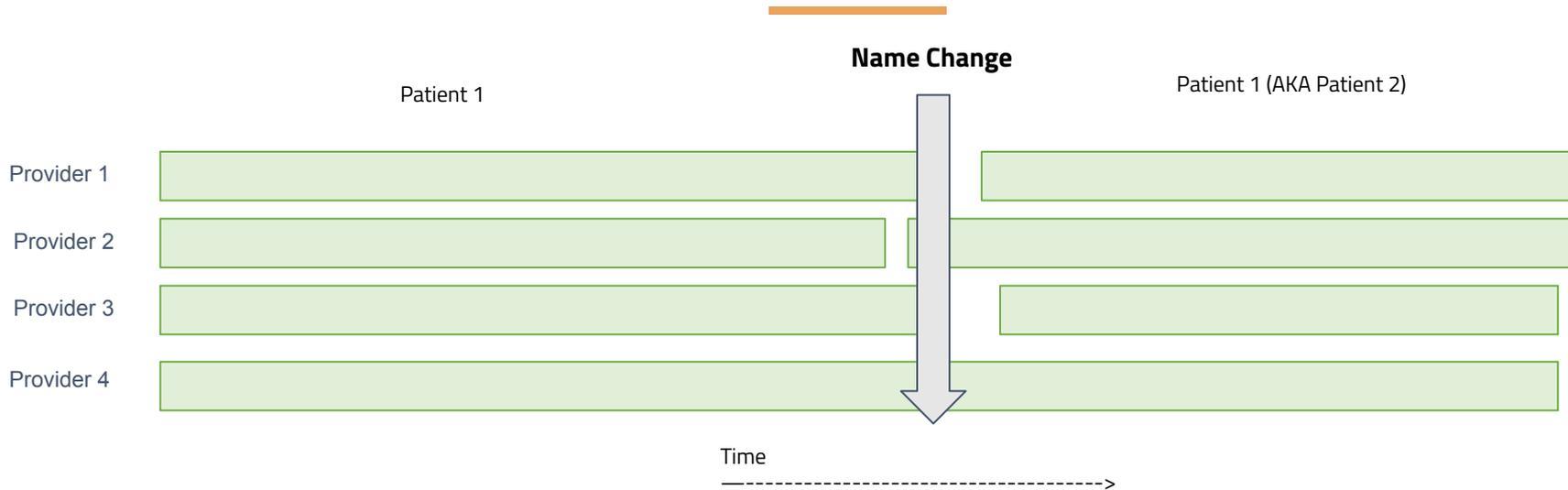
- # of claims
- # of claims overlapped with patient B
- # of chronic conditions similar to patient B
- # of chronic conditions exclusive to A

## Patient B:

- # of claims
- # of claims overlapped with patient A
- # of chronic conditions similar to patient A
- # of chronic conditions exclusive to B



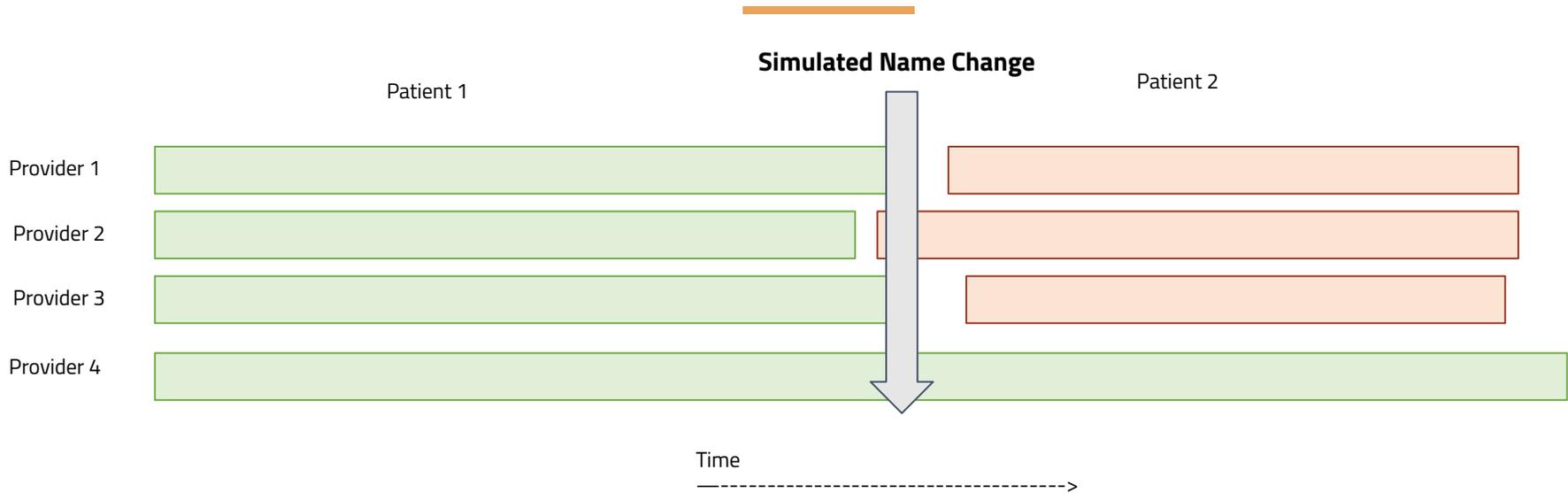
# Patient Dedupe Method #2: Synthetic Data



Simulate a “good” match by modifying closed claims to product synthetic training data

- Use a single patient and assign an arbitrary name change date
- Split the patient claims and create two new patients.

# Patient Dedupe Method #2: Synthetic Data

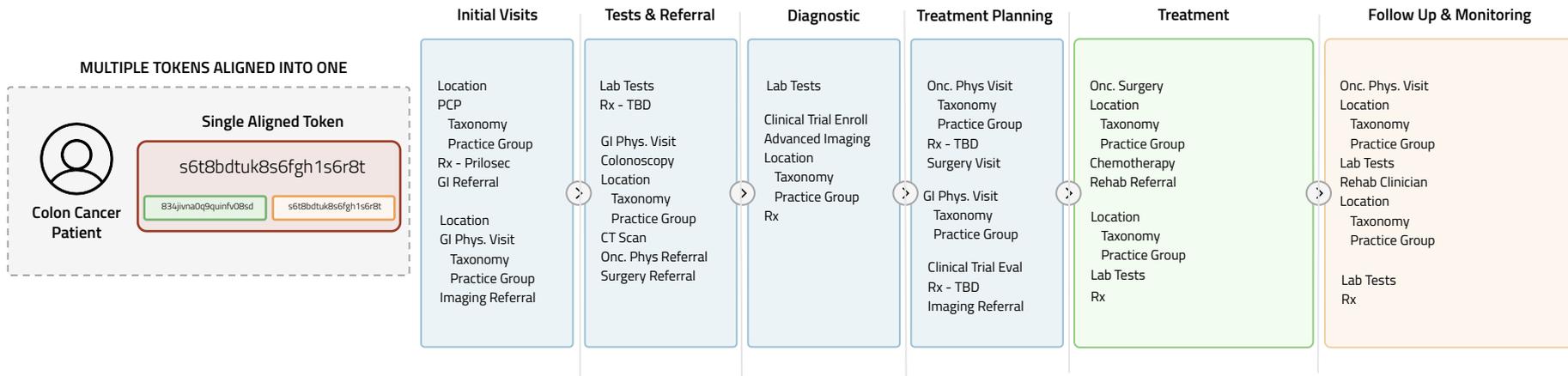


Simulate a "bad" match by modifying closed claims to product synthetic training data

- Use claims from two separate patients
- Truncate them before and after an arbitrary "name change date"

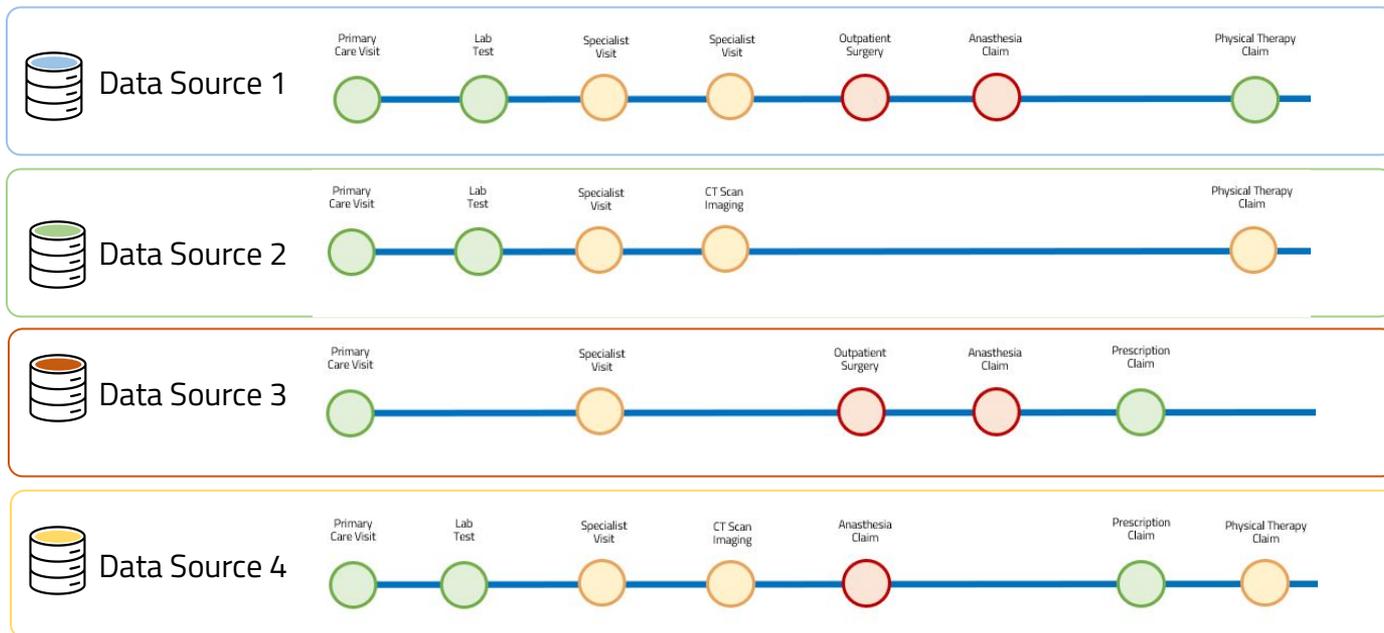
# Token Alignment

Token alignment leads to far more accurate event capture



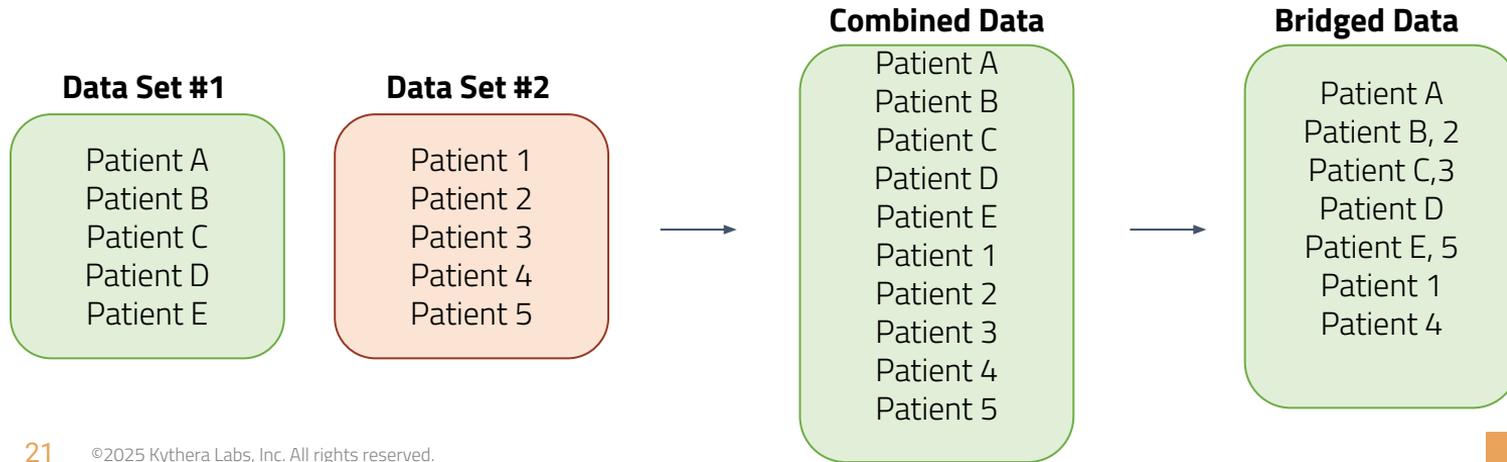
# Patient Journey Fragmented State: Multiple Sources

One Patient : Multiple Tokens : Multiple Data Sources

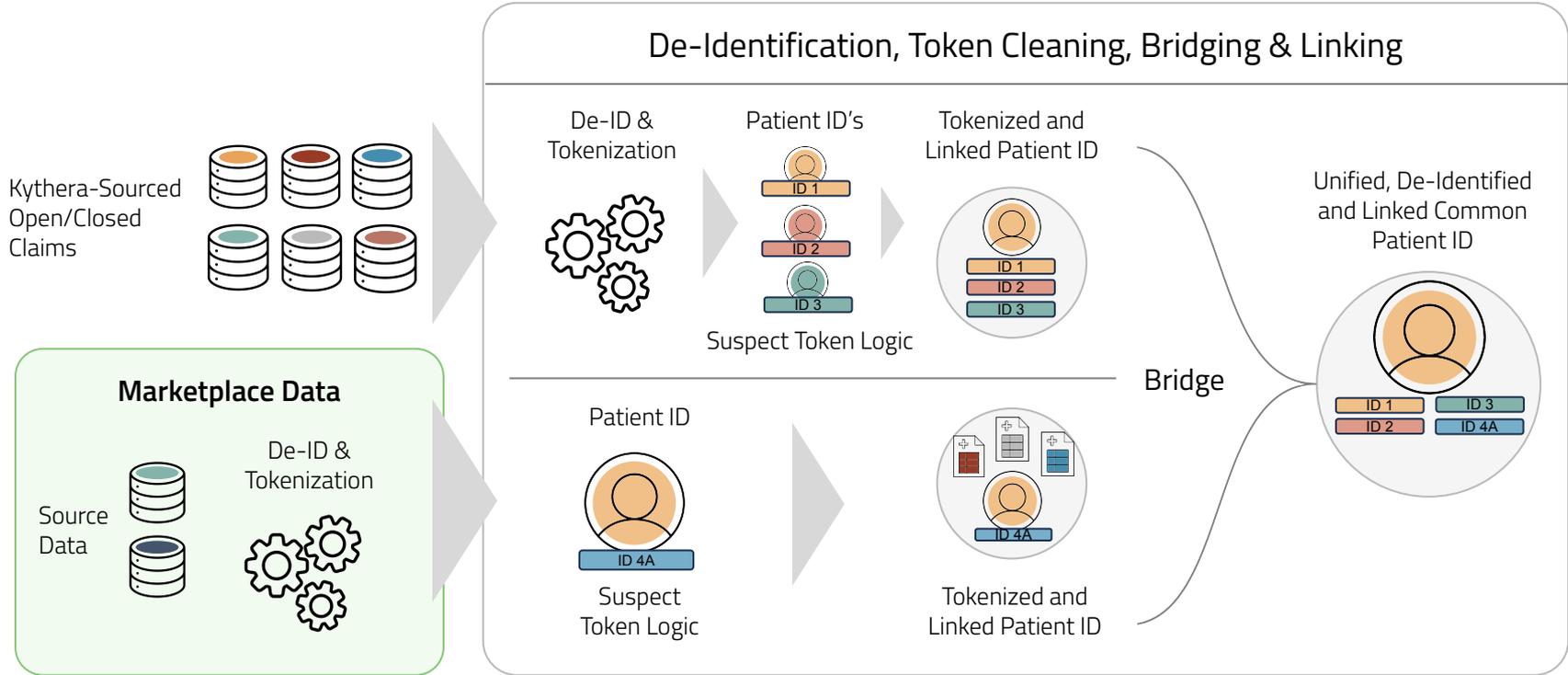


# Bridging Method

- First dedupe WITHIN each data set
- Treat separate data sets as one large data set
- Use token matching methods to tie tokens together



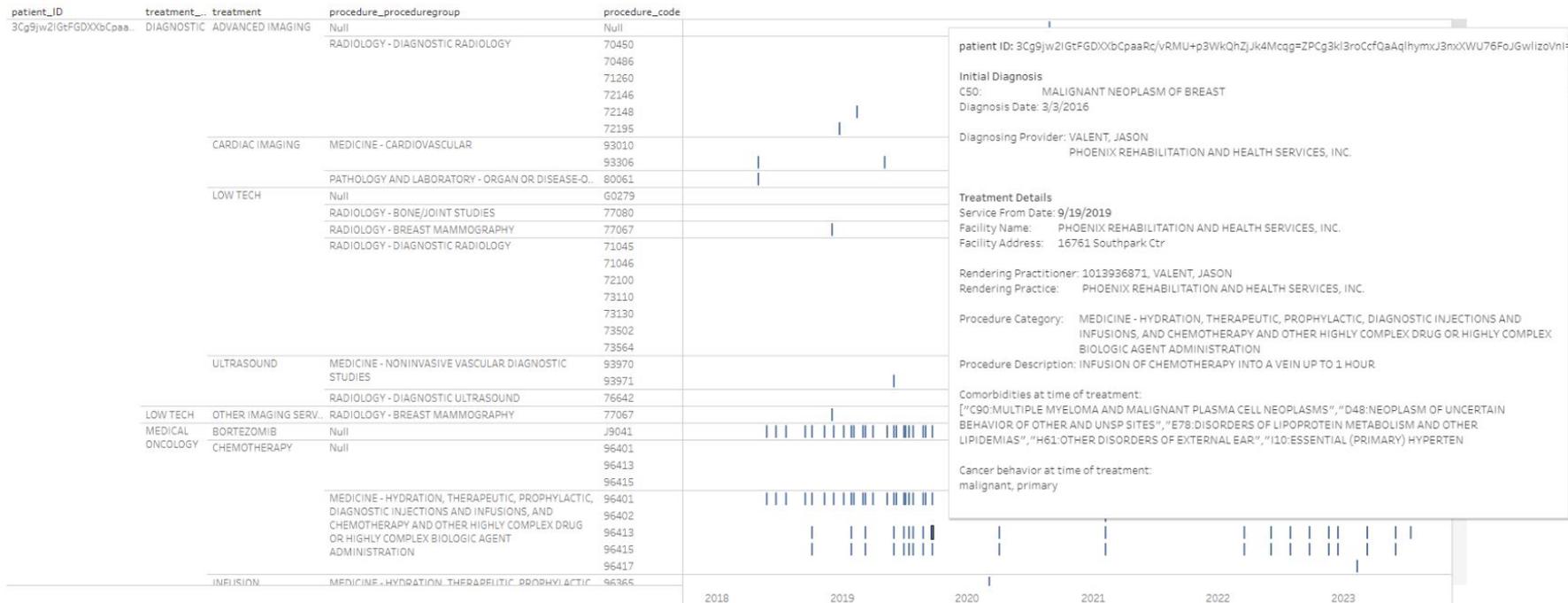
# Unifying Token Systems



# Patient Journey Example: Breast Cancer

Detail and visualize the patient journey

## Patient Clinical Pathway



# Understanding Unique Contribution of a Dataset

---

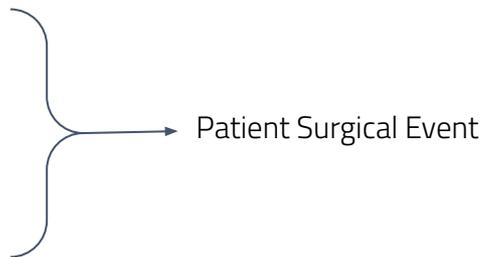
- “Overlap analysis” typically requested of vendors falls short
  - Prior to patient mastering overlaps are often misleading
  - Patient-only overlaps overestimate the value of a data combination
- Data overlaps best done at the “Patient Event” level
  - Better specificity to use case
  - Record-level data often redundant
  - Records can be of a different type
  - What matters is NEW INFORMATION

## Claims

Inst: Blood work  
Prof: Nursing care  
Prof: Anesthesia for surgery  
Prof: Surgeon  
Inst: Hospital Stay  
Inst: OR

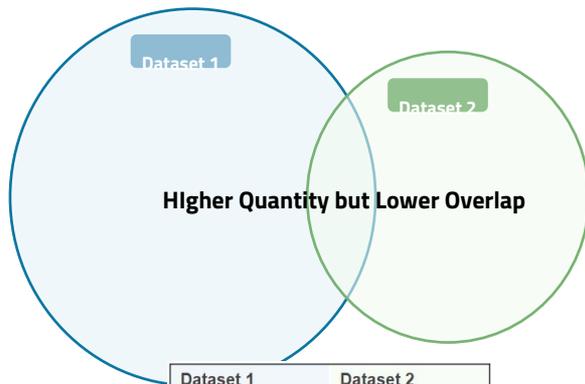
## EHR

Encounter: Surgery



# A Note About Overlaps

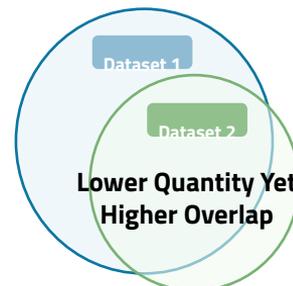
Before Patient Mastering



Dataset 1	Dataset 2
Token X	Token Z
Token Y	Token Q
Token R	

(Inflated Patient Count)

After Patient Mastering



Dataset 1	Dataset 2
Patient A	Patient A
Patient A	Patient B
Patient B	

(More Accurate Patient Count)

# Reducing Wasteful Data Spend through Source Assessment Intelligence: *Kythera Customer Case Study*

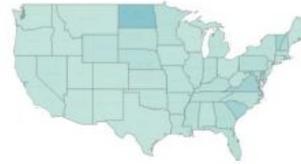
Unique Uplift of Surgical Event Claims by Claims Source and By State



> Source 1



> Source 2



> Source 3



> Source 4



> Source 5

Share Percent



0.0%

100.0%

# 30% reduction

in data spend was achieved  
by eliminating redundancy  
while optimizing their data  
investments to better  
support revenue-driving  
CPT codes.



Q & A

---

Contact:  
Stuart Head  
[stu@kytheralabs.com](mailto:stu@kytheralabs.com)

KYTERA



Thank You

---

KYTERA