

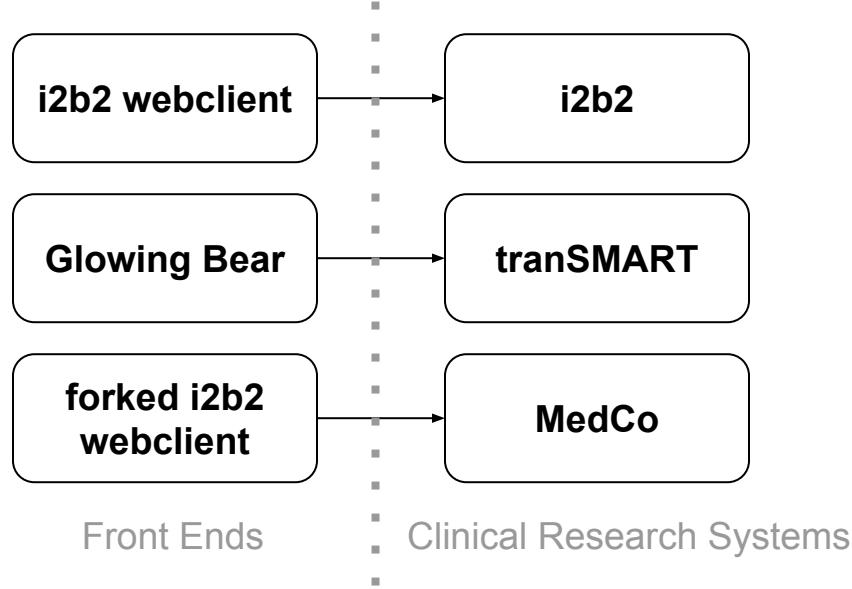
# Building a common and privacy-preserving front end for open-source clinical research platforms

*Presentation @ European i2b2 transSMART AUG, Geneva, 01/11/2018*

**Mickaël Misbach\*#, Ward Weistra#, Dr. Bo Gao#, Dr. Jean-Louis Raisaro\*,  
Dr. Juan Troncoso-Pastoriza\*, Prof. Jean-Pierre Hubaux\***

\*EPFL, #The Hyve

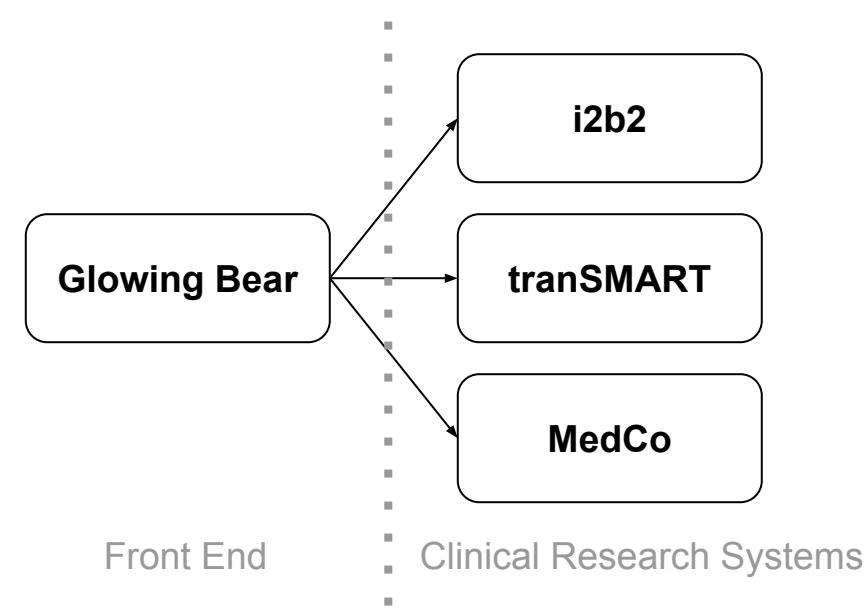
# Motivation



## Limitations:

- data fragmented in...
  1. location
  2. technical solution
- sensitive data not easily shared

# Objectives



**Goals:** Extend Glowing Bear capabilities

- integrate additional clinical research platform: i2b2
- enable privacy-preserving cohort exploration: MedCo

**Why:** Enable scientists to access more data from a common interface by...

- expanding compatibility of data source backends
- accessing sensitive data that would otherwise be difficult to share

**How:** Add a layer of interoperability on top of Glowing Bear to support i2b2 and MedCo

# Building Blocks

# Common User Interface: Glowing Bear



The screenshot shows the Glowing Bear user interface. At the top, there's a navigation bar with tabs for "Data Selection" and "Analysis". The "Data Selection" tab is active. On the far right of the header, it says "0.0.1-SNAPSHOT" with a help icon and a back arrow. Below the header, there's a search bar labeled "Specify query name" and a "Save query" button. The main area is divided into sections:

- Current Data Selection:** Buttons for "... subjects" and "... observations" with a "Clear all" button.
- Ontology:** A sidebar with a "filter" and "clear" button, and a tree view of ontology categories:
  - Vital Signs
    - Heart Rate (4)
  - Public Studies
    - CATEGORICAL\_VALUES (3)
    - CLINICAL\_TRIAL (3)
      - Demography
      - Vital Signs
        - Heart Rate (3)
    - CLINICAL\_TRIAL\_HIGHDIM (3)
    - EHR (3)
    - EHR\_HIGHDIM (3)
    - MIX\_HD (3)
    - Oracle\_1000\_Patient (1,200)
    - RNASEQ\_TRANSCRIPT (3)
    - SHARED\_CONCEPTS\_STUDY
    - SHARED\_CONCEPTS\_STUDY
    - SHARED\_HD\_CONCEPTS\_STU
    - SHARED\_HD\_CONCEPTS\_STU
    - TUMOR\_NORMAL\_SAMPLES
  - Projects
  - Private Studies
- Step 1: Define subjects**: Shows "8 / 1,246 subjects (1%) , 26 / 120,188 observations (0%)". An "Update" button is next to it.
- Inclusion criteria:** 8 subjects included. This section contains two complex search criteria:
  - Concept:** Gender (\Projects\Survey 1\Demographics) with values Female (5).  
with more options
  - Study:** SURVEY1

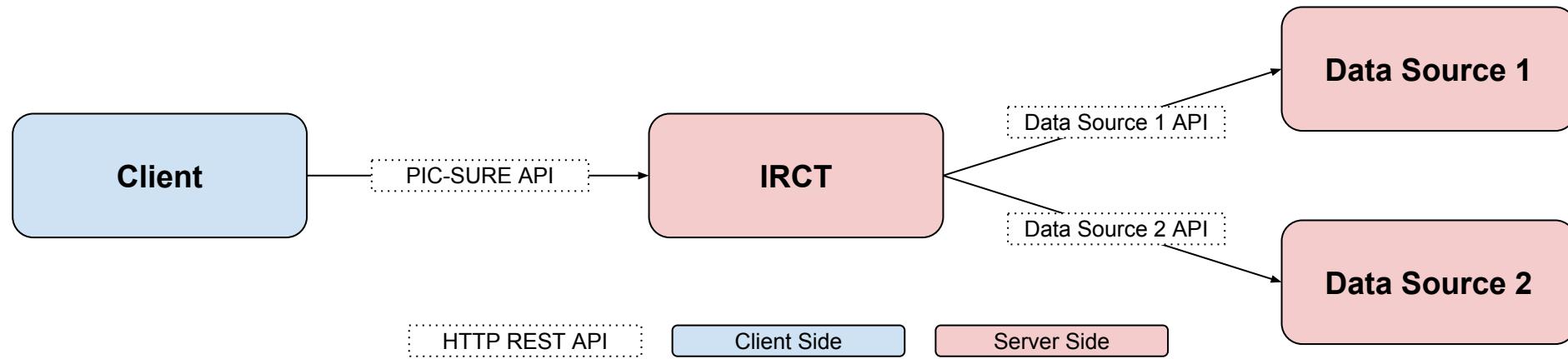
An "Import Criteria" button is located at the top right of this section.

**or**

  - Concept:** Heart Rate (\Public Studies\CLINICAL\_TRIAL\Vital Signs) with value between 20 and max.  
with more options

# Common Query Language: PIC-SURE API[1]

- PIC-SURE provides a common API to query any kind of *data sources*
- At the technical level, data semantic is not covered



PIC-SURE: Patient-centered Information Commons: Standardized Unification of Research Elements

IRCT: Inter-Resource Communication Tool

HMS-DBMI: Harvard Medical School - Department of Biomedical Informatics

[1]: Alex AT Bui, John Darrell Van Horn, NIH BD2K Centers Consortium, et al. “Envisioning the future of big data biomedicine”. In: Journal of biomedical informatics 69 (2017), pp. 115–117.

# i2b2[1] & tranSMART[2]

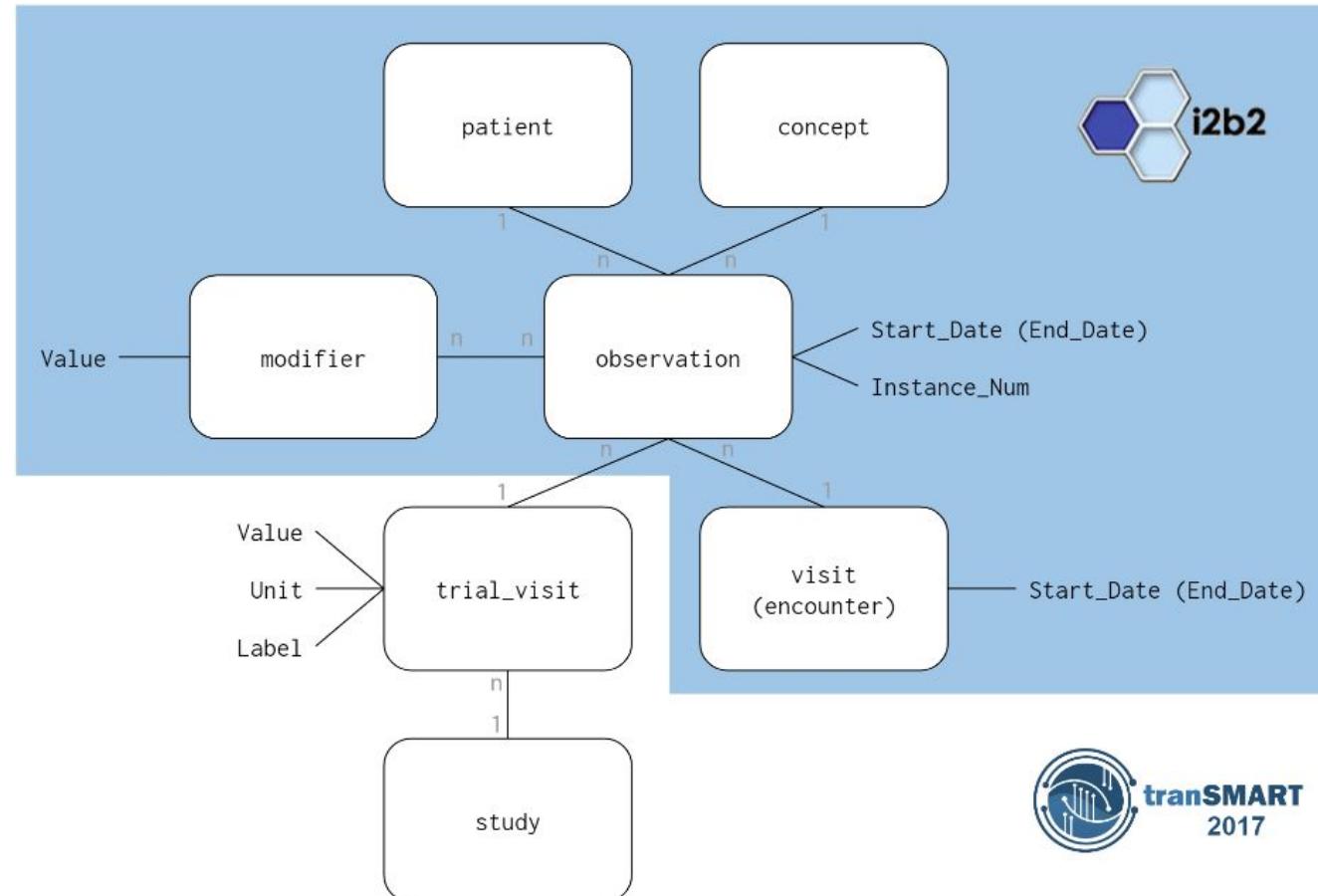


## i2b2

- cohort exploration
- database: star schema
- preference from hospitals

## tranSMART

- i2b2 features and...
- advanced cohort exploration features
- advanced data export
- preference from pharma
- study-based



[1]: Shawn N Murphy et al. "Serving the enterprise and beyond with informatics for integrating biology and the bedside (i2b2)". In: Journal of the American Medical Informatics Association 17.2 (2010), pp. 124–130.

[2]: Elisabeth Scheufele et al. "tranSMART: an open source knowledge management and high content data analytics platform". In: AMIA Summits on Translational Science Proceedings 2014 (2014), p. 96.

# MedCo[1]



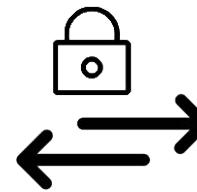
Distributed cohort exploration



Secure storage outsourcing



Trust decentralization



End-to-end data protection



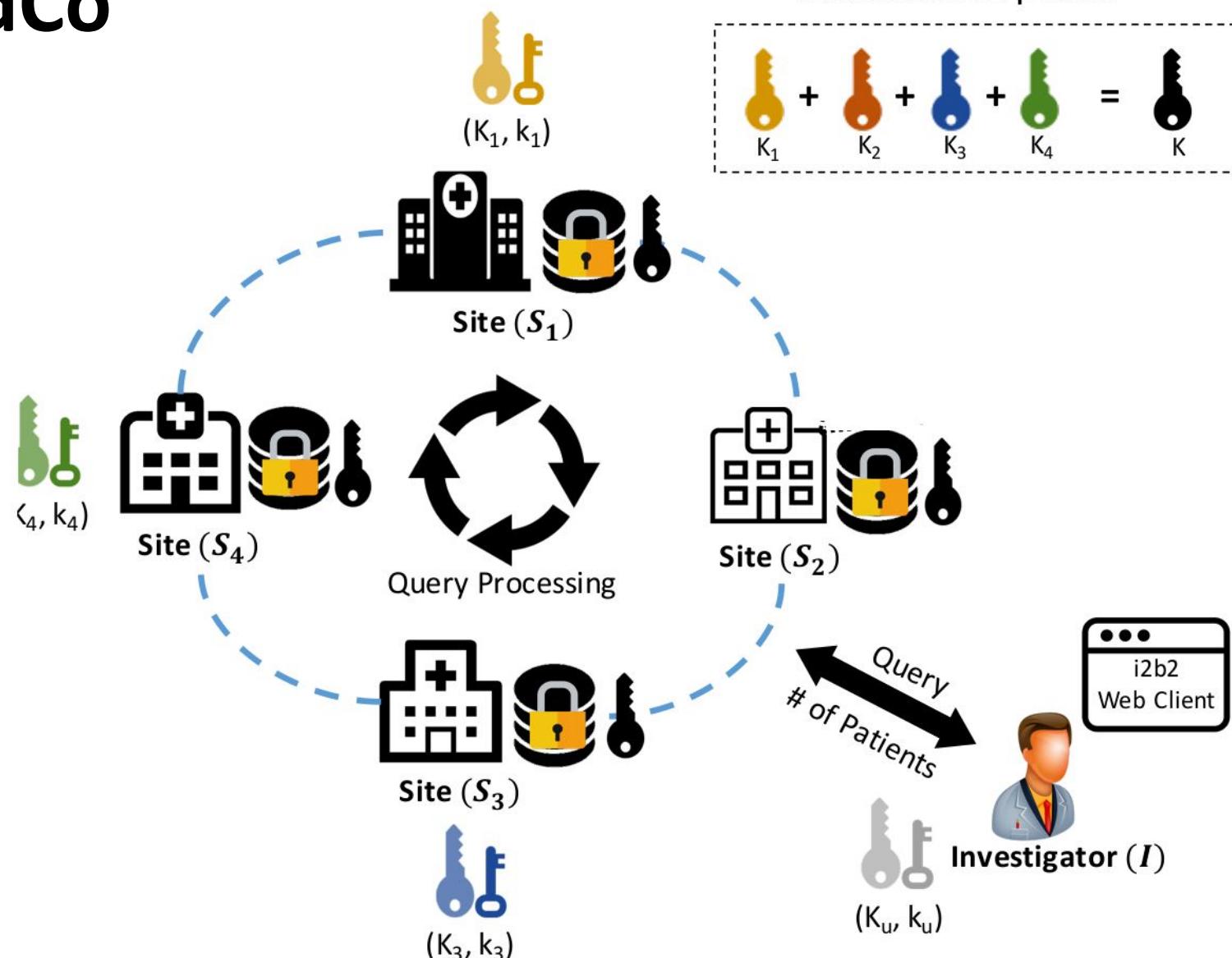
Unlinkability

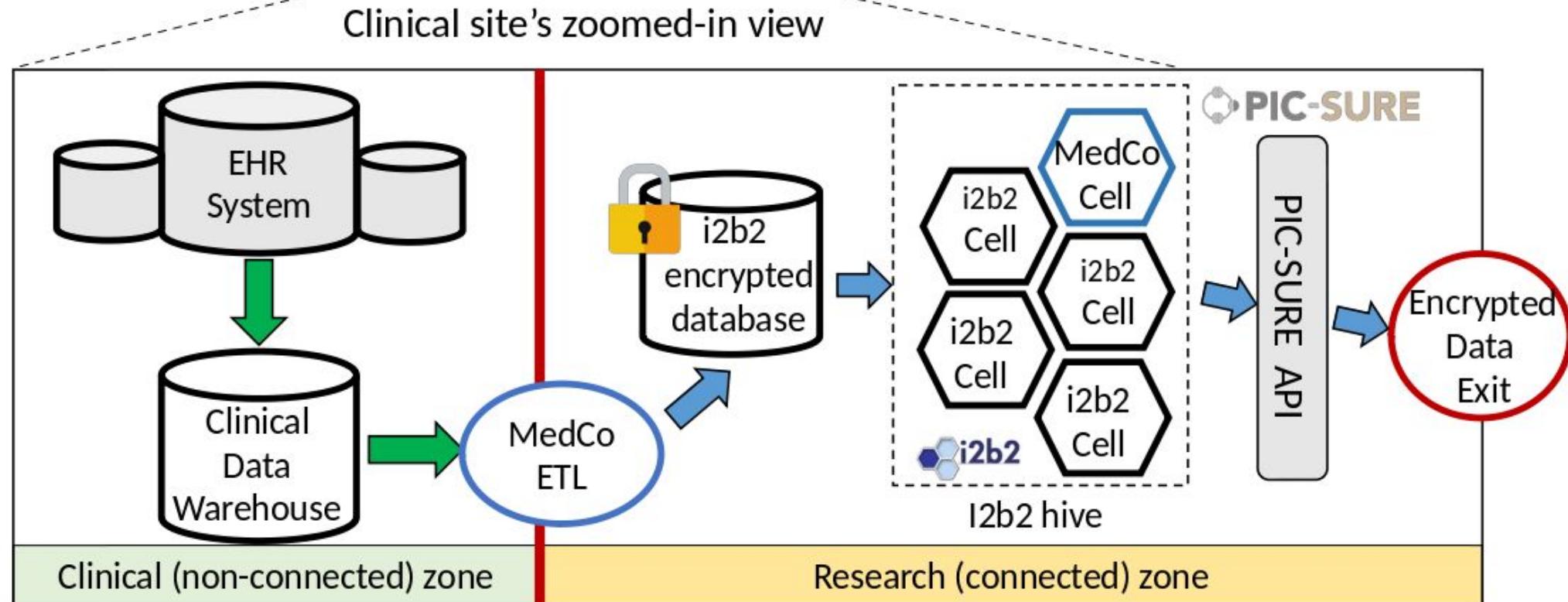


Differential privacy

⇒ MedCo provides technical means to share data that would otherwise not be shared

[1]: J. L. Raisaro et al. "MedCo: Enabling Secure and Privacy-Preserving Exploration of Distributed Clinical and Genomic Data". In: IEEE/ACM Transactions on Computational Biology and Bioinformatics (2018), pp. 1–1. issn: 1545-5963. doi: 10.1109/TCBB.2018.2854776.





ETL: Extraction Transformation  
(Encryption) and Loading tool



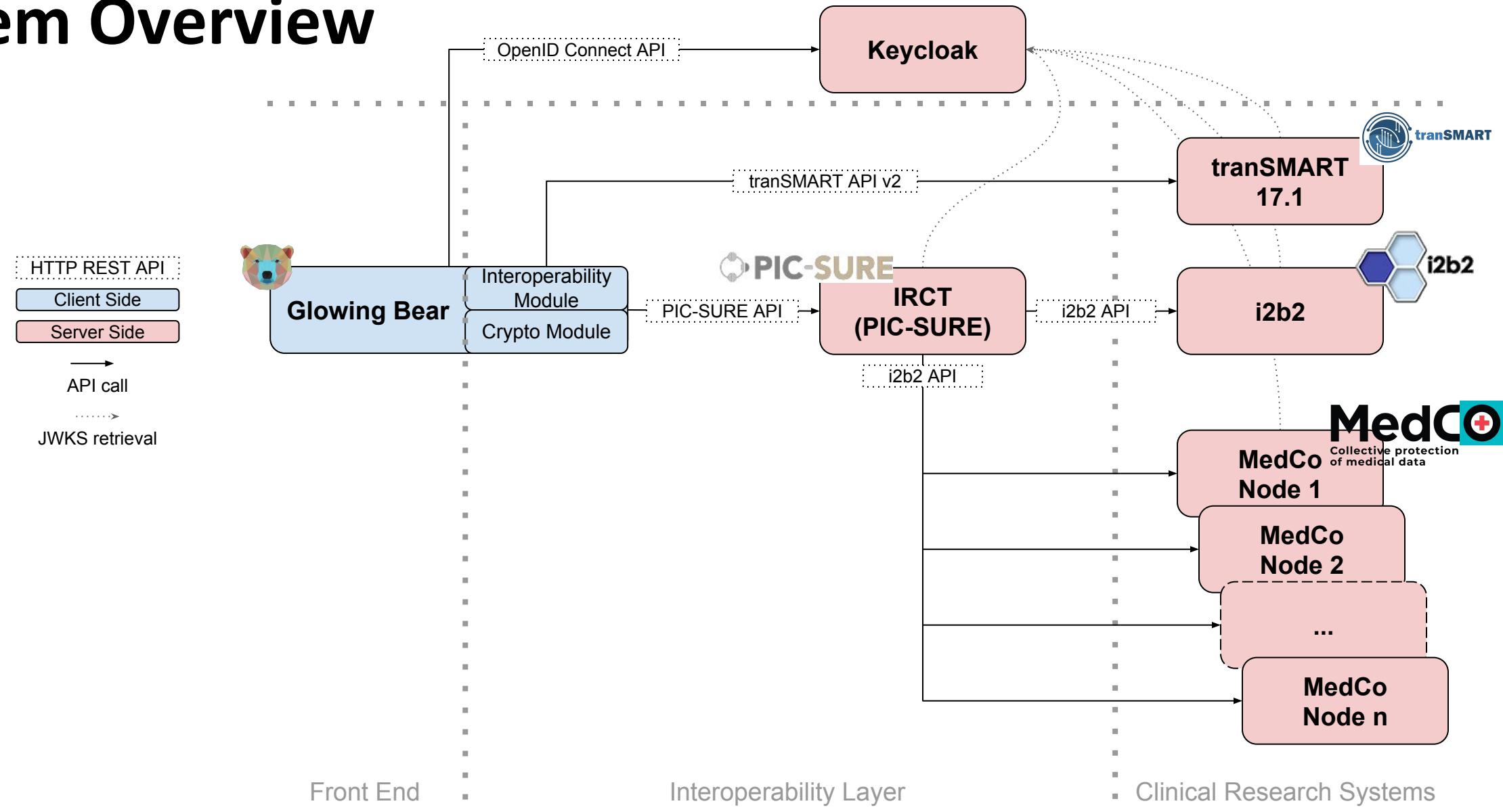
Clear data



Encrypted data

# System Overview

# System Overview



# Demo

# Conclusion

- offer a modern front end that allows cohort exploration: **Glowing Bear** ✓
- compatible with **tranSMART** (v17.1) and **i2b2** ✓
- compatible with **MedCo** ✓
- be extensible for **future support of additional platforms** ✓
- technical considerations:
  1. easy to deploy ✓
  2. not degrade user experience in existing systems ✓
  3. enforce secure authentication ✓
  4. open-source ✓
  5. practical runtime ✓

# Future Work

- Authentication: using a distributed ledger, to avoid single point of failure (Keycloak)  
(in progress)
- Support additional systems through the PIC-SURE API  
(in progress: HAIL[1], framework for exploration and analysis of genomic data)
- More features in UI for i2b2 and MedCo  
e.g. data export, analysis, query saving, etc.
- PIC-SURE 2.0 upgrade

# Context

## **Project is a collaboration between EPFL & The Hyve**

- The Hyve: wanted Glowing Bear to support i2b2
- EPFL: wanted new front end for MedCo
- Within the framework of the Swiss DPPH (Data Protection in Personalized Health)

**talk with chair guy for presentation  
part of work being sponsor by dpph**

# User Interface before: i2b2 webclient

i2b2 Query & Analysis Tool    Project: Genomics – WES    User: i2b2 Demo User    Find Patients | Analysis Tools | Message Log | Help | Logout

**Query Tool**

Query Name:

Temporal Constraint: Treat all groups independently

Group 1			Group 2			Group 3		
Dates	Occurs > 0x	Exclude	Dates	Occurs > 0x	Exclude	Dates	Occurs > 0x	Exclude
Treat Independently			Treat Independently			Treat Independently		
HLA-DQB1 levels < 0			SNV/SNP [stop_gained] SNV/SNP [non_synonymous]			Gene Symbol [Contains: HLA-DQB1]]		

Run Query   Clear   Print Query   3 Groups   New Group

**Query Status**

Number of patients  
**28**  
For Query "HLA-D-SNV/S-SNV/S@03:50:12"

**Navigate Terms**   **Find Terms**

- Sequence Ontology Variants
  - complex substitution
  - copy number variation
  - deletion
  - indel
  - insertion
  - inversion
  - MNP
- Reference genome
- SNV/SNP
  - 3 prime UTR variant
  - 5 prime UTR variant
  - Alternate Allele
  - Chromosome
  - dbSNP RS id
  - downstream gene variant
  - exon variant
  - frameshift\_variant
  - HGNC Gene Symbol
  - inframe\_variant
  - intergenic variant
  - intron variant
  - non\_synonymous
  - PolyPhen2 prediction
  - PolyPhen2 score
  - Reference Allele
  - Sequence End (base pair)
  - Sequence Start (base pair)
  - splicing variant
  - stop\_gained
  - stop\_lost
  - synonymous
  - upstream gene variant
  - Zygoticity
  - point\_mutation

# MedCo



**Website:** *medco.epfl.ch*

## Roadmap

- MedCo with PIC-SURE and Glowing Bear: Nov. 2018
- MedCo with MedChains (blockchain-based authentication and access control): April 2019
- Skype-like pull model: June 2019

# OpenID Connect[1] / Keycloak

## OpenID Connect

- authentication and authorization protocol
- token-based
- allow any kind of services to delegate identity and access management

## Keycloak

- identity provider
- implements OpenID Connect server

## Why?

- many different components running on different systems, with different authentication mechanisms
- need to have common authentication

# System Design Overview



# Glowing Bear[1]

- User interface for cohort exploration
- And more advanced analytics
- Open-source web application
- Originally for tranSMART v17.1 only

# PIC-SURE API

- Query terms of data source exposed through a *tree*
- Each data source declares its query format
- Data source has freedom of implementing anything as long as it fits the interface
- API has a SQL-like format:

```
"where": [ {  
    "field": {  
        "pui": "/resource/study/Age/",  
        "dataType": "INTEGER"  
    },  
    "predicate": "CONSTRAIN_VALUE",  
    "fields": {  
        "OPERATOR": "GT",  
        "CONSTRAINT": "20"  
    }  
}]
```

where part: constraints on data

field: path and type of query term queried  
(obtained from tree)

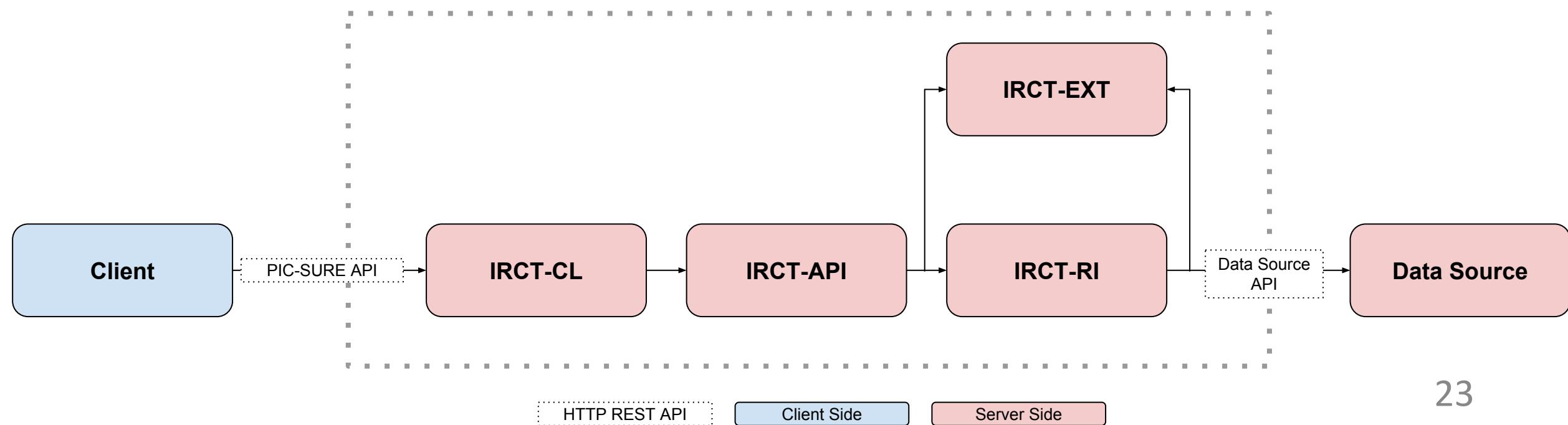
*predicate* used on query term

fields: additional input to predicate

# IRCT

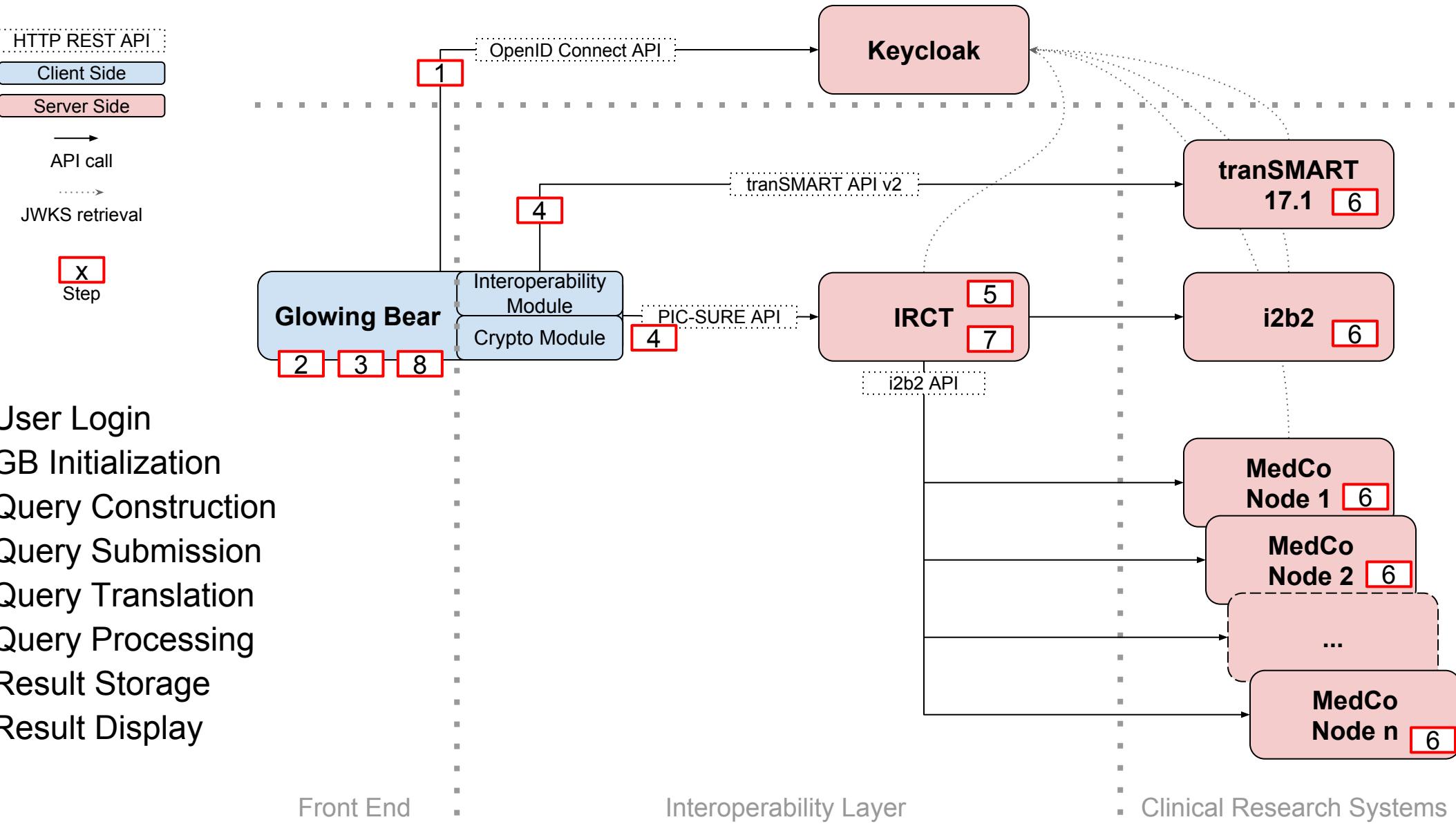
- IRCT is the official implementation of the PIC-SURE API
- 4 components:
  - IRCT-CL (REST service)
  - IRCT-API (core library)
  - IRCT-EXT (external hooks library)
  - IRCT-RI (data sources connectors)

**IRCT**: Inter-Resource Communication Tool  
**CL**: Communication Layer  
**API**: Application Programming Interface  
**RI**: Resource Interface  
**EXT**: EXTension

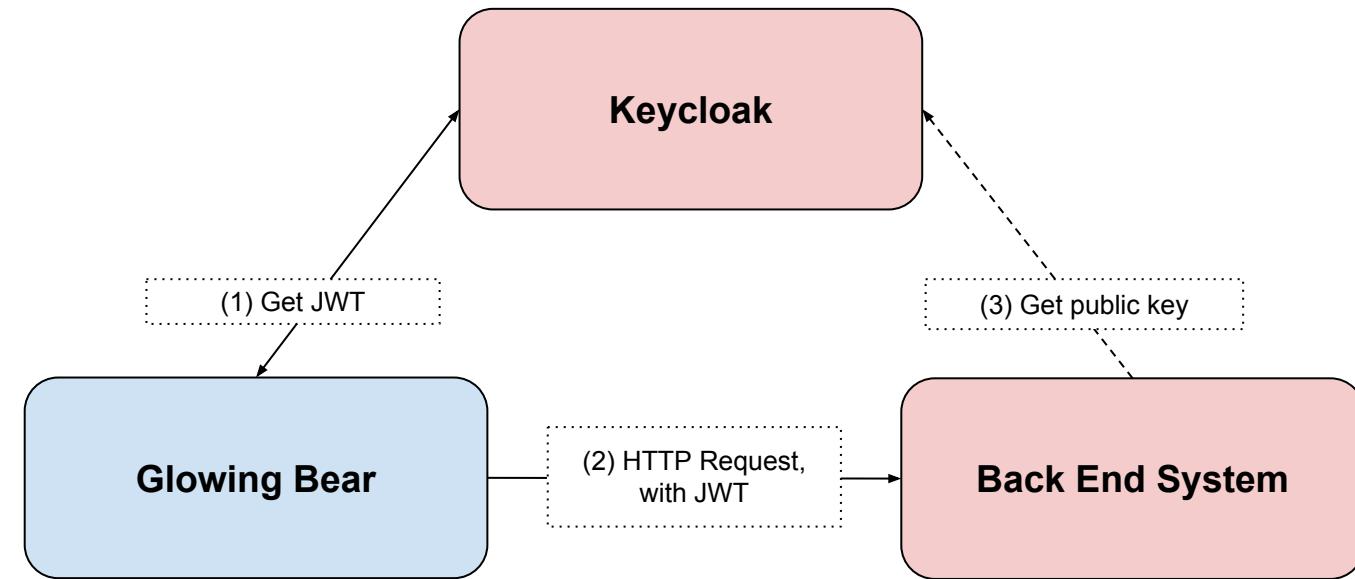


# Detailed Workflows

# Detailed Query Workflow



# OpenID Connect: Stateless Authentication



JWT: JSON Web Token

### Legend

-----> Cached Request

HTTP Request

Server Side

Client Side

# OpenID Connect: JSON Web Token

## JWT format

→ 3 distinct base64-encoded values

- JSON header: metadata

```
{  
  "alg": "RS256",  
  "typ": "JWT",  
  "kid": "eTFrdyrNxXLNHI7p0Ywybc7z1SBHTEcqWcMTybtqvQY"  
}
```

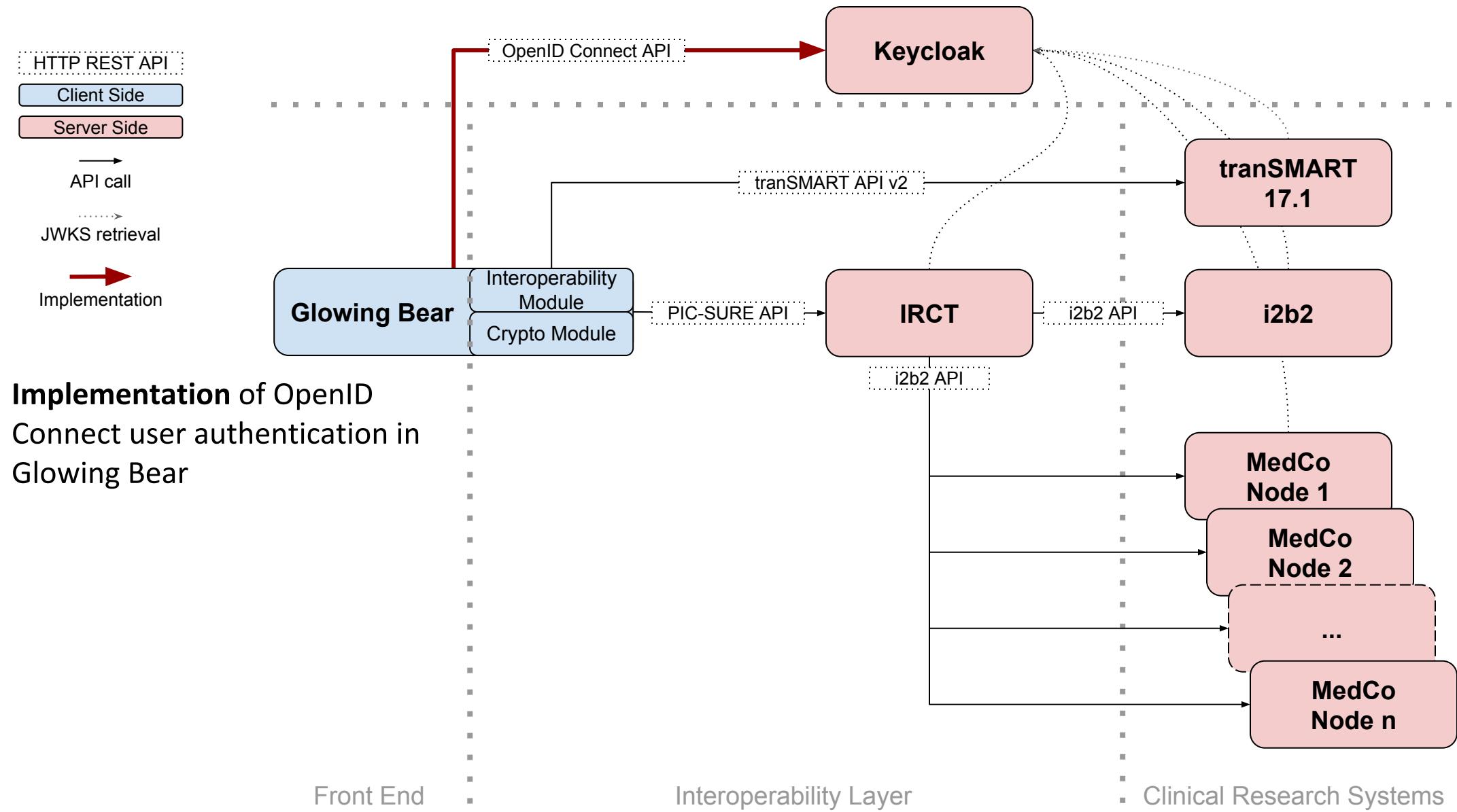
- JSON payload: identity, authorizations, validity, ...

```
{  
  "exp": 1523454086,  
  "iat": 1523453186,  
  "iss": "http://localhost:8081/auth/realms/master",  
  "aud": "glowing-bear",  
  "nonce": "N0.28573339803406971523453198656",  
  "resource_access": {  
    "account": {  
      "roles": [  
        "role1",  
        "role2"  
      ]  
    } },  
  "preferred_username": "test",  
  "email": "test@test.com"  
}
```

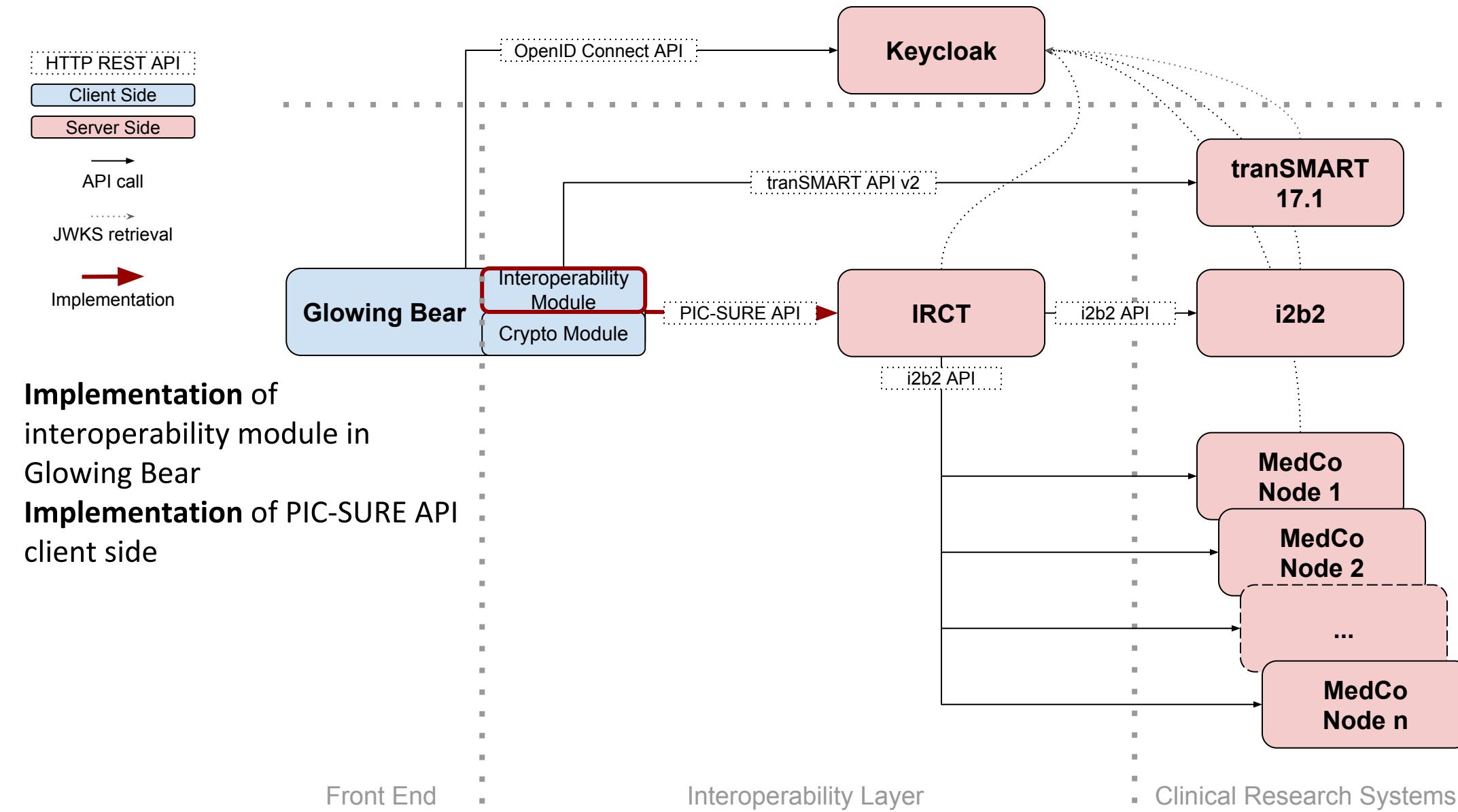
- Binary signature

# Implementation

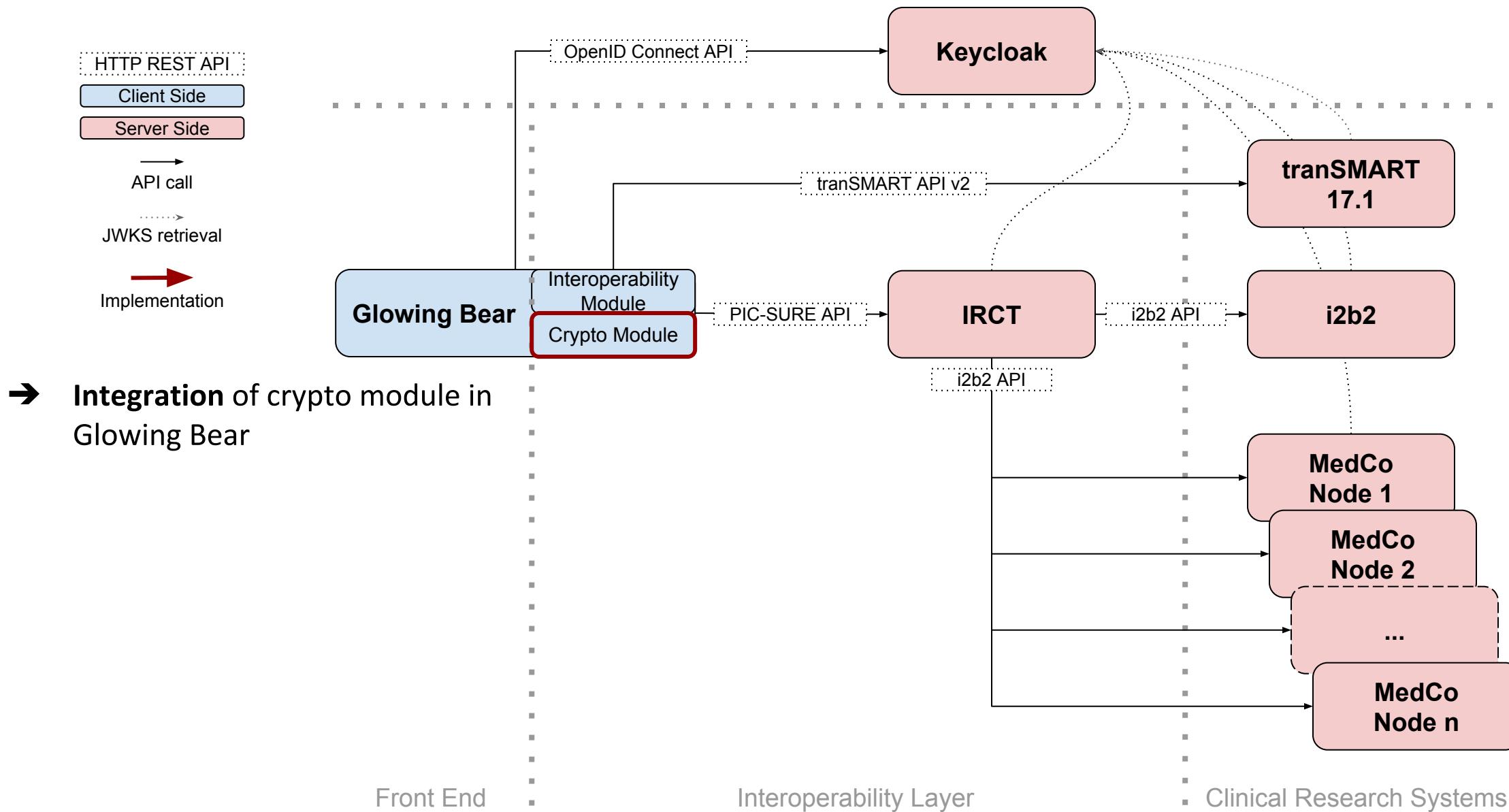
# Glowing Bear: OpenID Connect



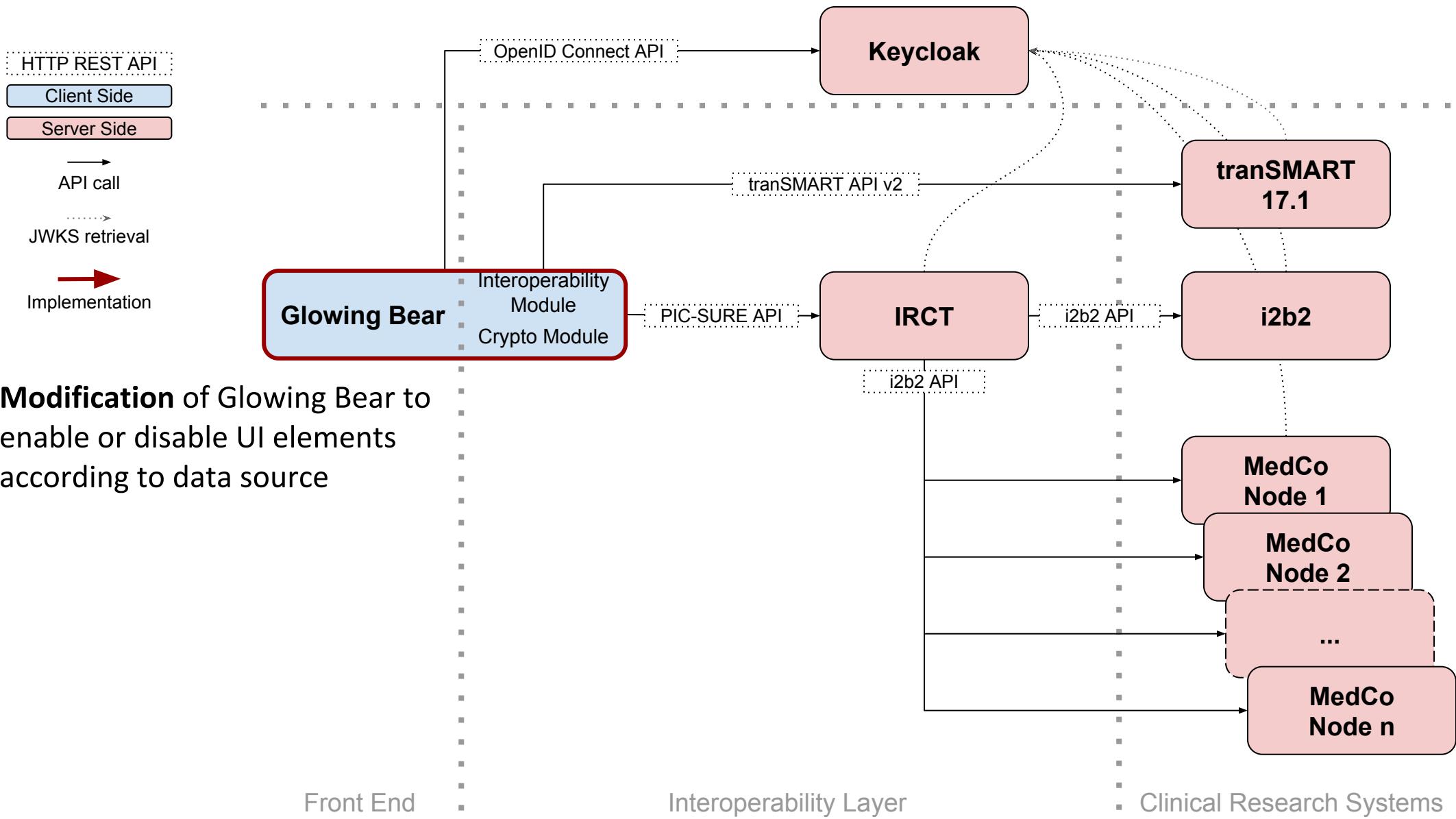
# Glowing Bear: Interoperability Module



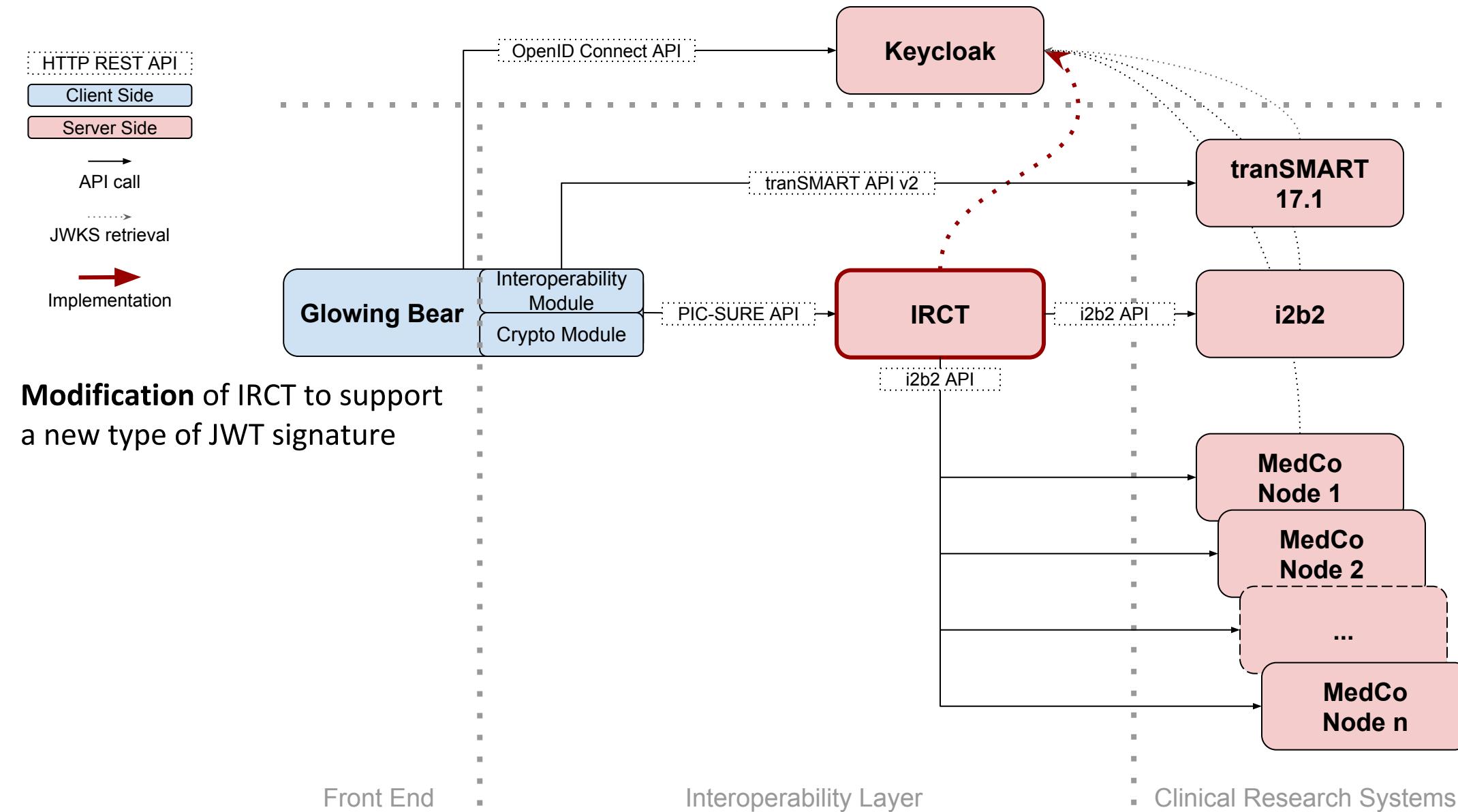
# Glowing Bear: Crypto Module



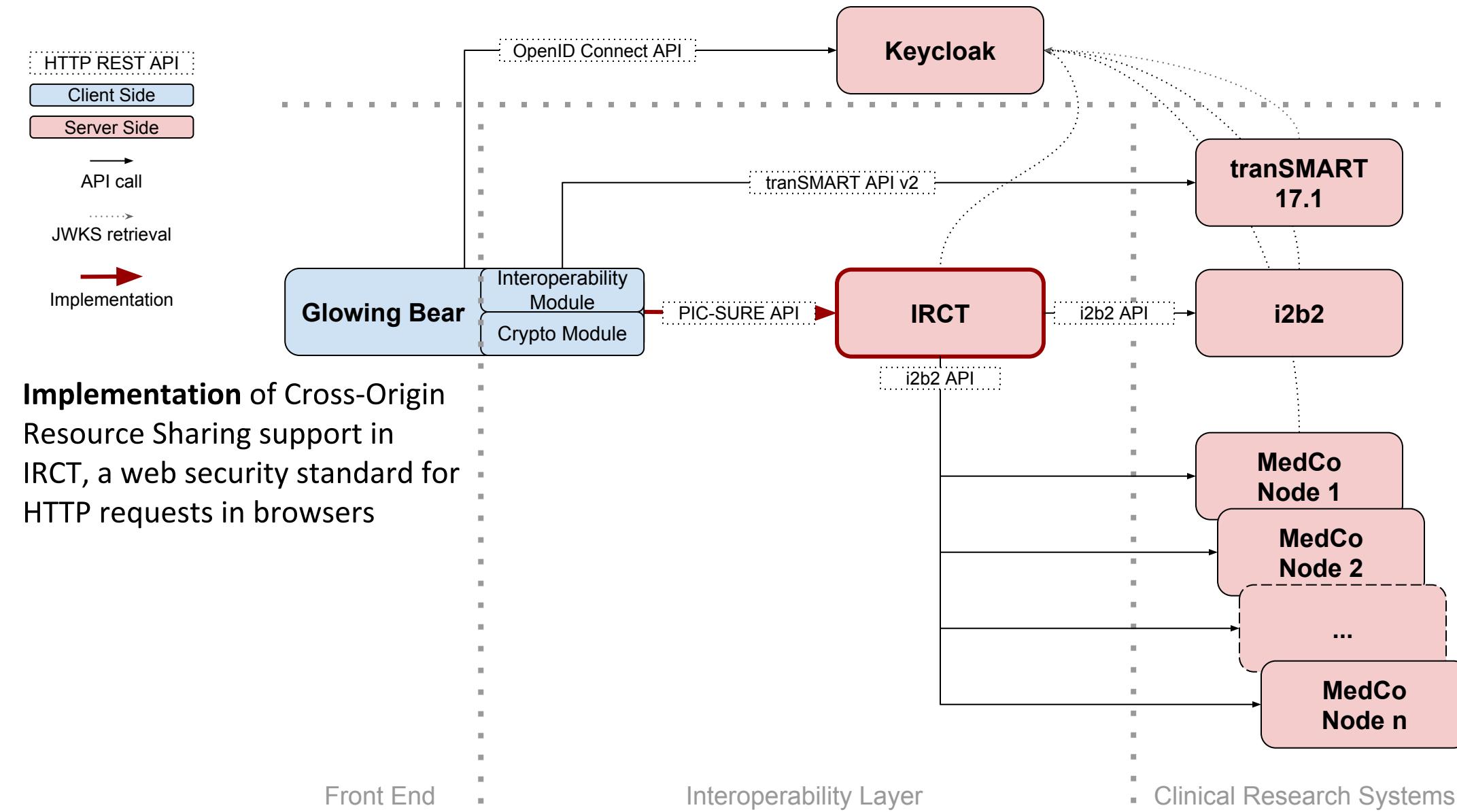
# Glowing Bear: User Interface



# IRCT: OpenID Connect

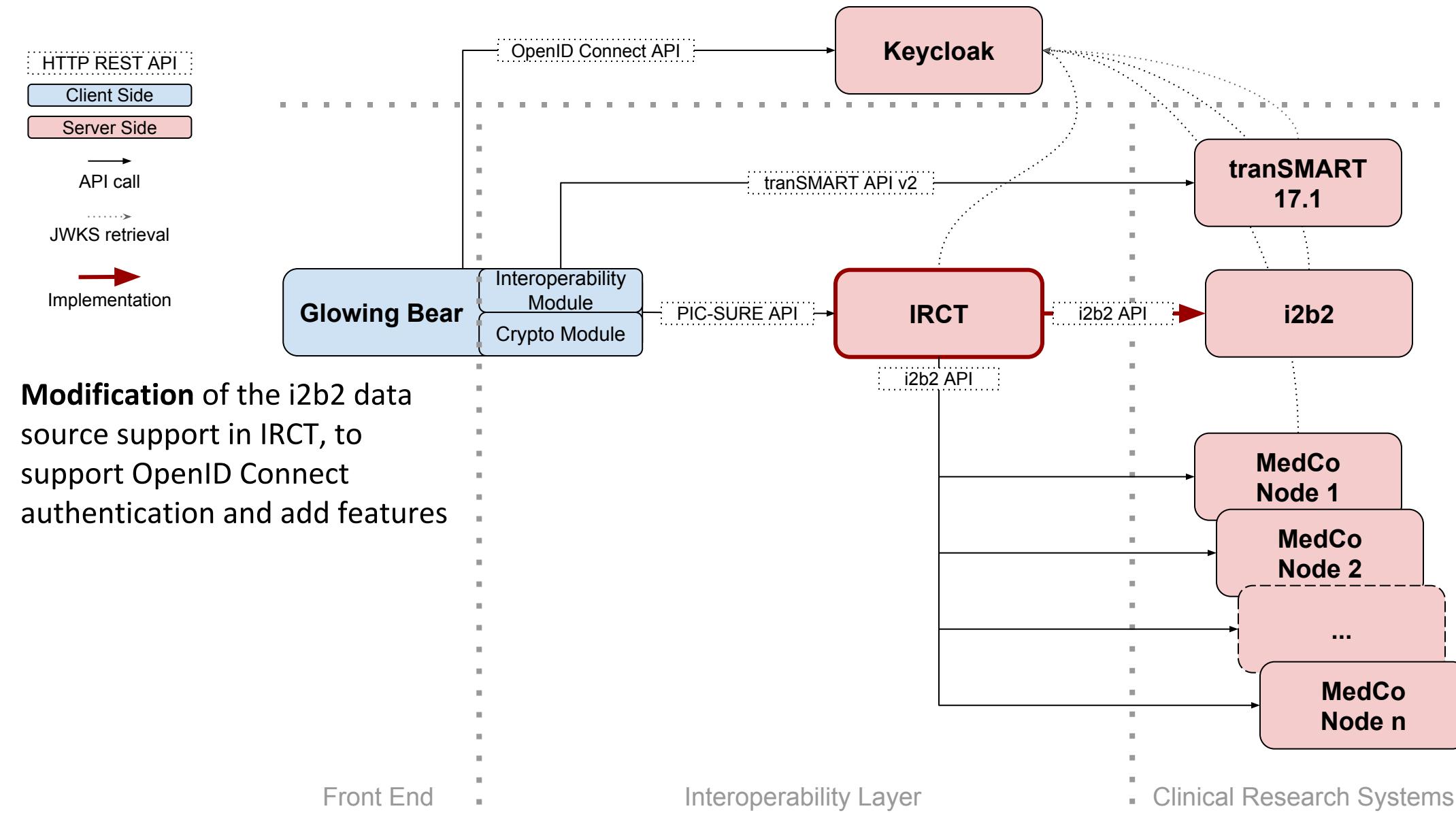


# IRCT: Cross-Origin Resource Sharing (CORS)

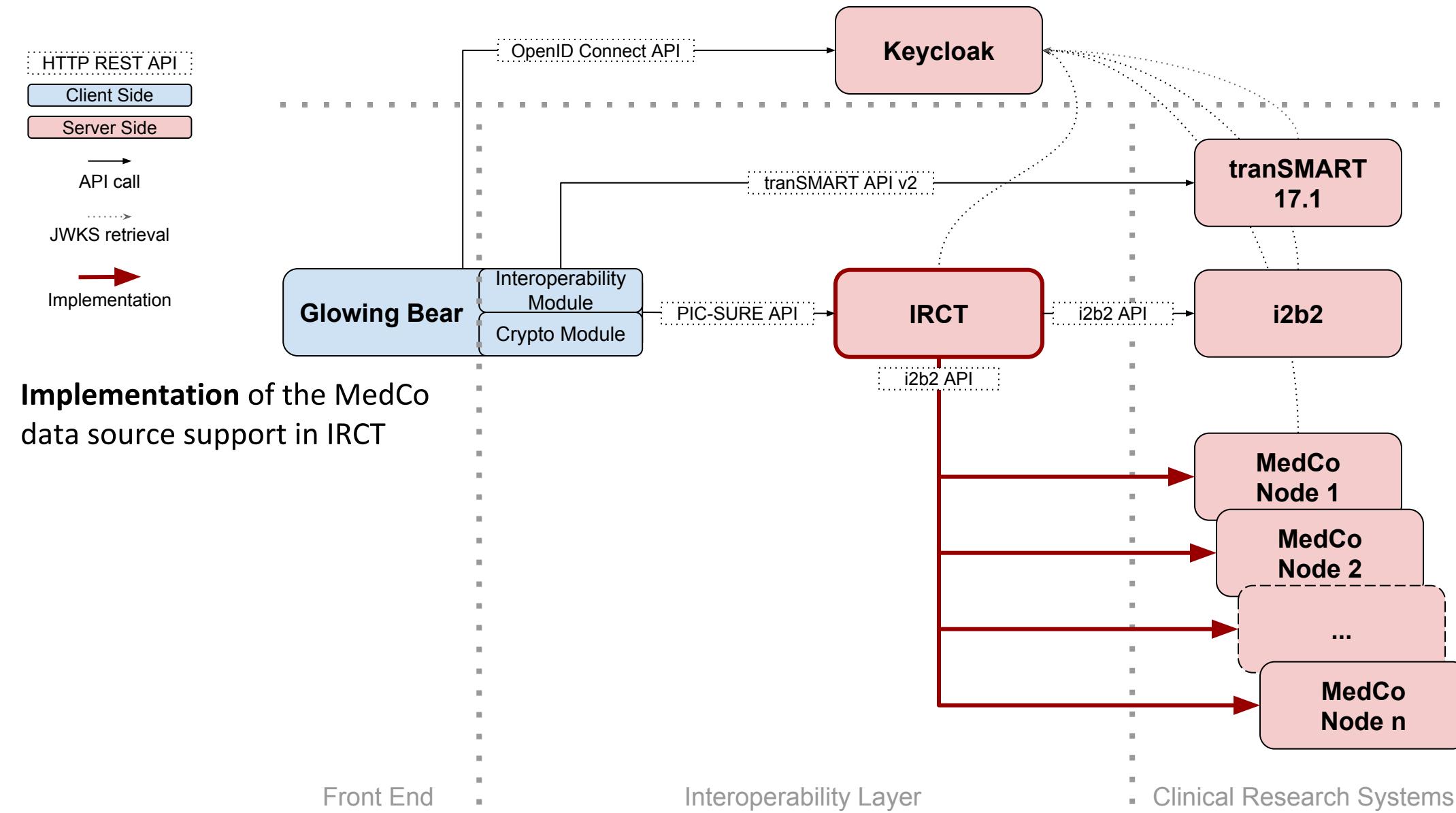


→ **Implementation** of Cross-Origin Resource Sharing support in IRCT, a web security standard for HTTP requests in browsers

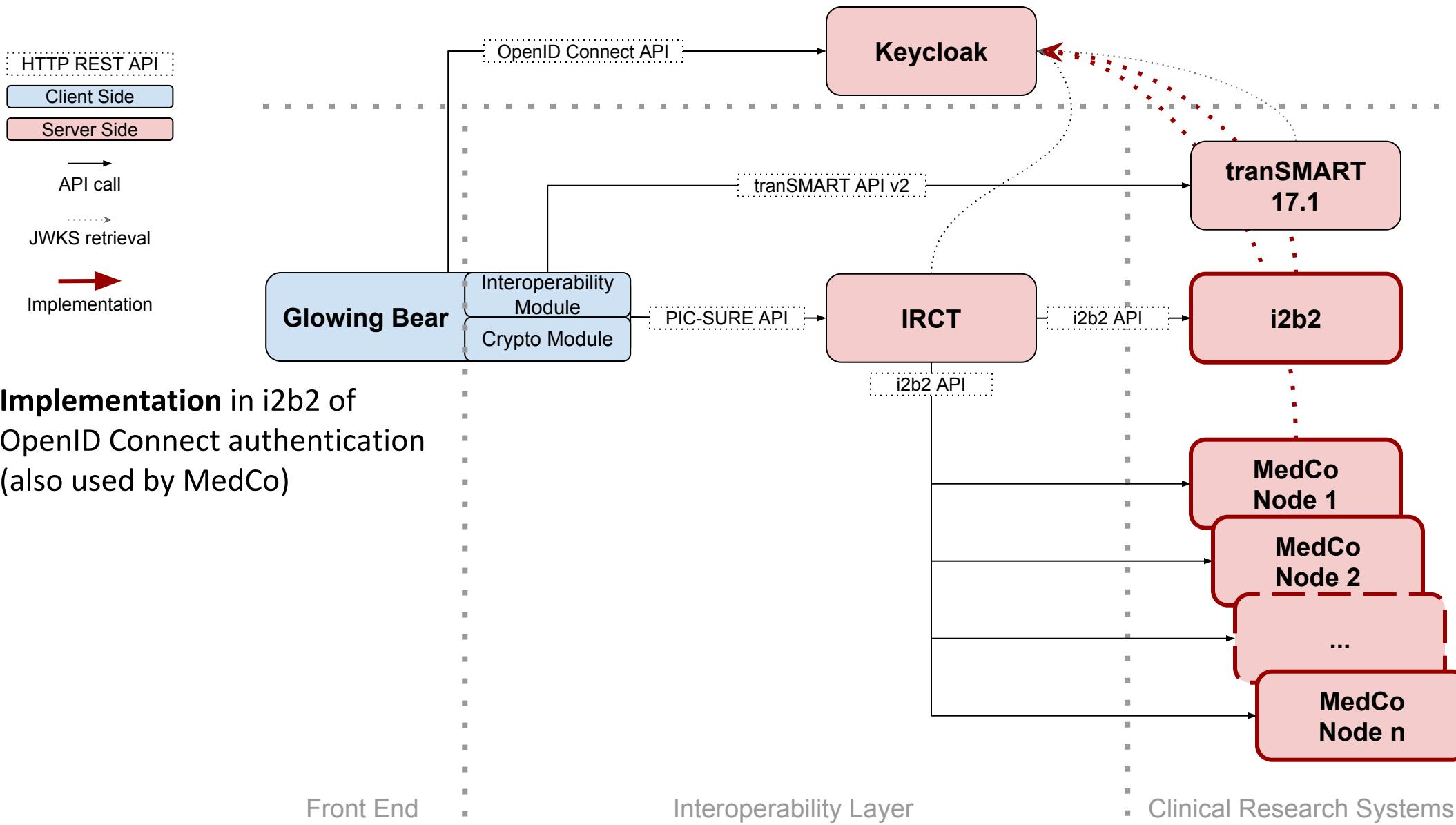
# IRCT: i2b2 Data Source



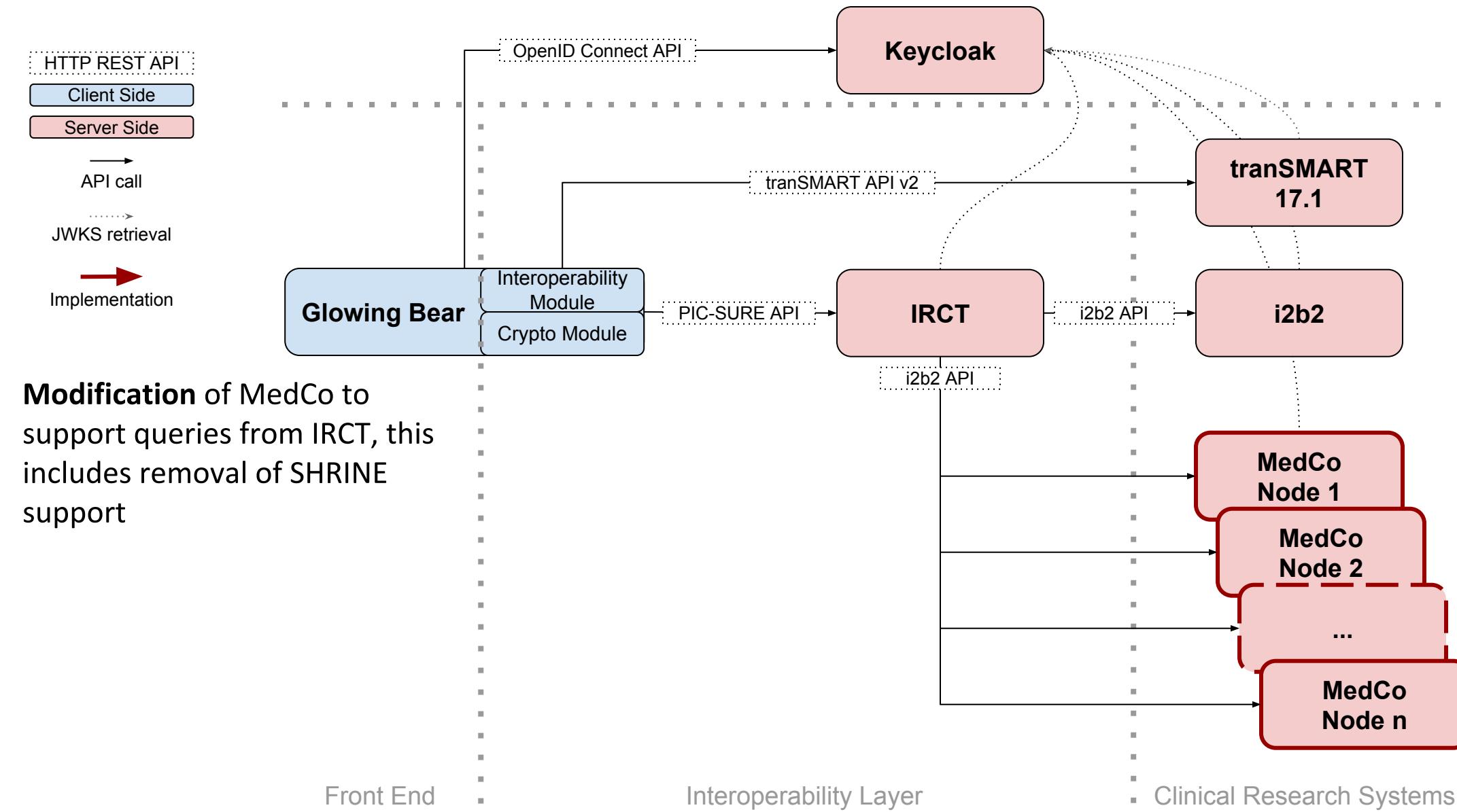
# IRCT: MedCo Data Source



# i2b2: OpenID Connect



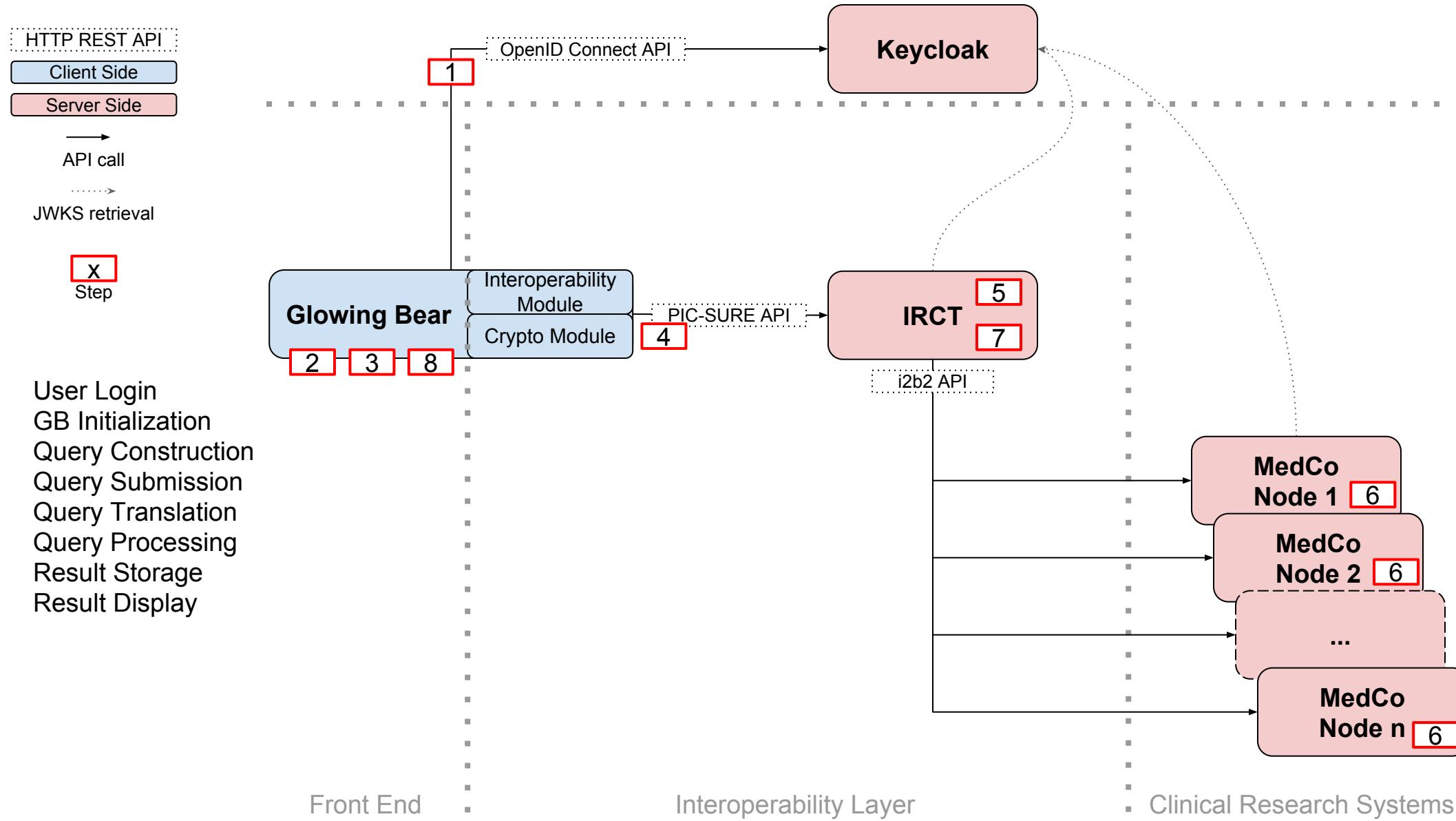
# MedCo: IRCT Querying



# Conclusion

# Backups

# Query Workflow: MedCo



# Query Workflow: MedCo

## Step 2: Glowing Bear Initialization

- Public key of the collective authority loaded (collective authority = all the MedCo nodes):  
 $Pk_c$
- Pair of public/private key of user randomly generated  
 $Pk_u / pk_u$

## Step 3: Query Construction

- In the tree, each query term has an a corresponding integer  
 $q_v$
- It is encrypted before submission with the public key of the collective authority  
 $ENC_{PkC}[q_v]$

## Step 5: Query Translation

- Query in PIC-SURE API format translated to i2b2 API format
- Broadcasted to all the MedCo nodes at the same time

# Query Workflow: MedCo

## Step 6: Query Processing (at each node, using distributed protocols involving all nodes)

- A. Encrypted query terms are “tagged”: the encryption is switched from probabilistic to deterministic (DDT = Distributed Deterministic Tagging)

$$\text{ENC}_{\text{PkC}}[q_v] \rightarrow \text{DDT}_{S_i}[q_v]$$

- B. The query is submitted to normal i2b2 (in the database, the  $\text{DDT}_{S_i}[q_v]$  are stored)
- C. i2b2 answers with a patient set, but it contains dummy patients

# Query Workflow: MedCo

## Step 6: Query Processing (distributed protocols involving all the MedCo nodes)

D. Each patient has a “dummy flag”: an encrypted 0 or 1

→ we fetch the flags of patients from the set

$$\text{ENC}_{\text{PkC}}[f_j]$$

E. The sum of the flags is our real count, we compute it homomorphically:

$$\text{ENC}_{\text{PkC}}[R_i]$$

F. The result is encrypted with the key of the collective authority, the corresponding private key does not exist

→ we change the key of the encryption to the one of the user

$$\text{ENC}_{\text{PkU}}[R_i]$$

# Query Workflow: MedCo

## Step 8: Result Display

- The encrypted results are fetched by Glowing Bear
- They are decrypted using the private key of the user

$R_i$