# BIPMed: Bioinformatics challenges and solutions

Benilton S Carvalho, Ph.D.

Institute of Mathematics, Statistics and Scientific Computing

University of Campinas



BIPMed **Brazilian Initiative on PRECISION MEDICINE**

# Genomic Databases

- LOVD interface offers a lot of extra annotation;

- Loading data to this system requires a lot of computational resources;

- DB performance may be sub-optimum on the client-side;

- Testing transition to column-based DBMS;



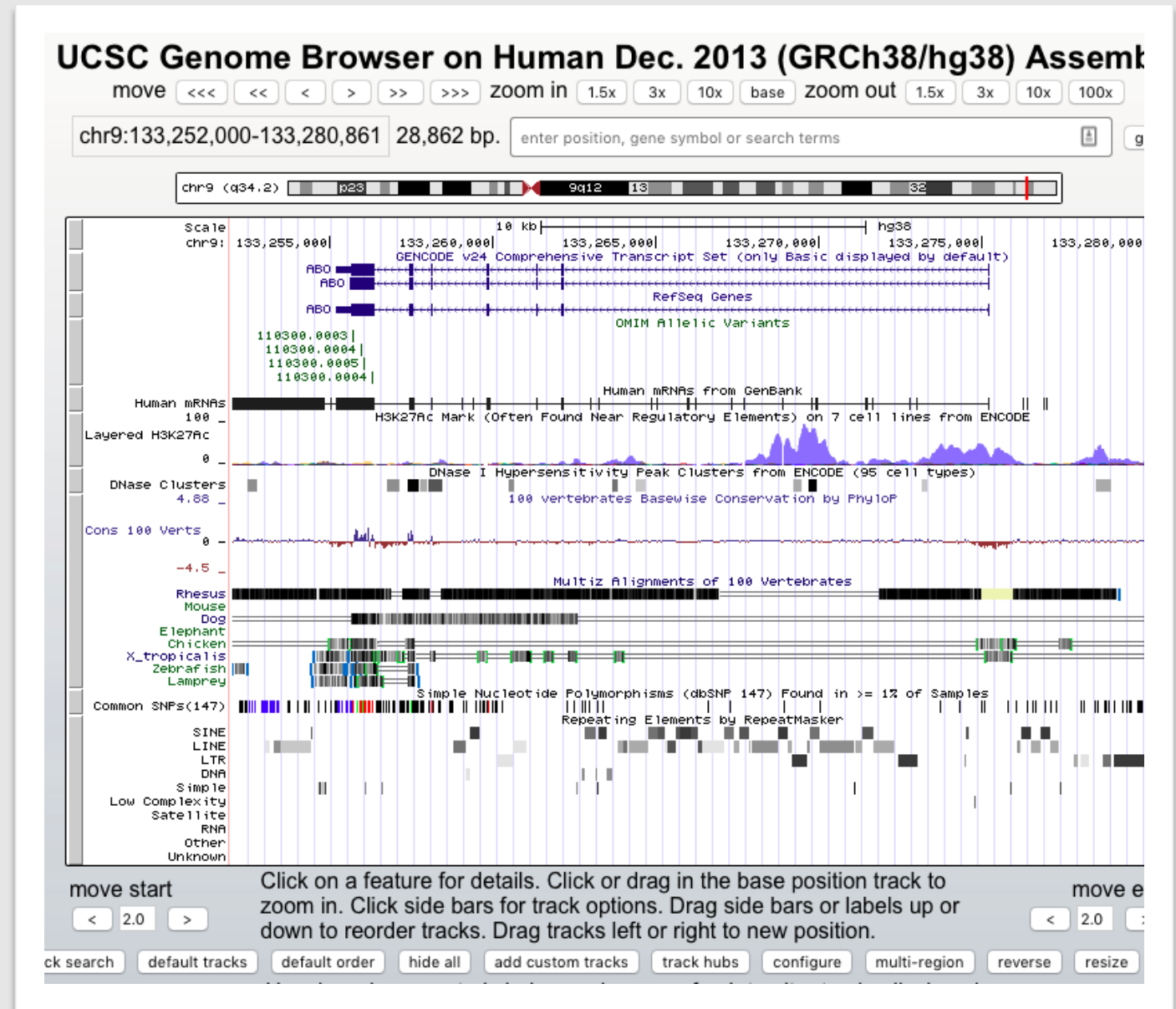**Brazilian Initiative on Precision Medicine**

contact@bipmed.org

| | Variants | Individuals | Diseases | Screenings | Submit | Documentation |

Showing entries 1 - 100.

« First | ‹ Prev | **1** | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | ... | Next › | Last »

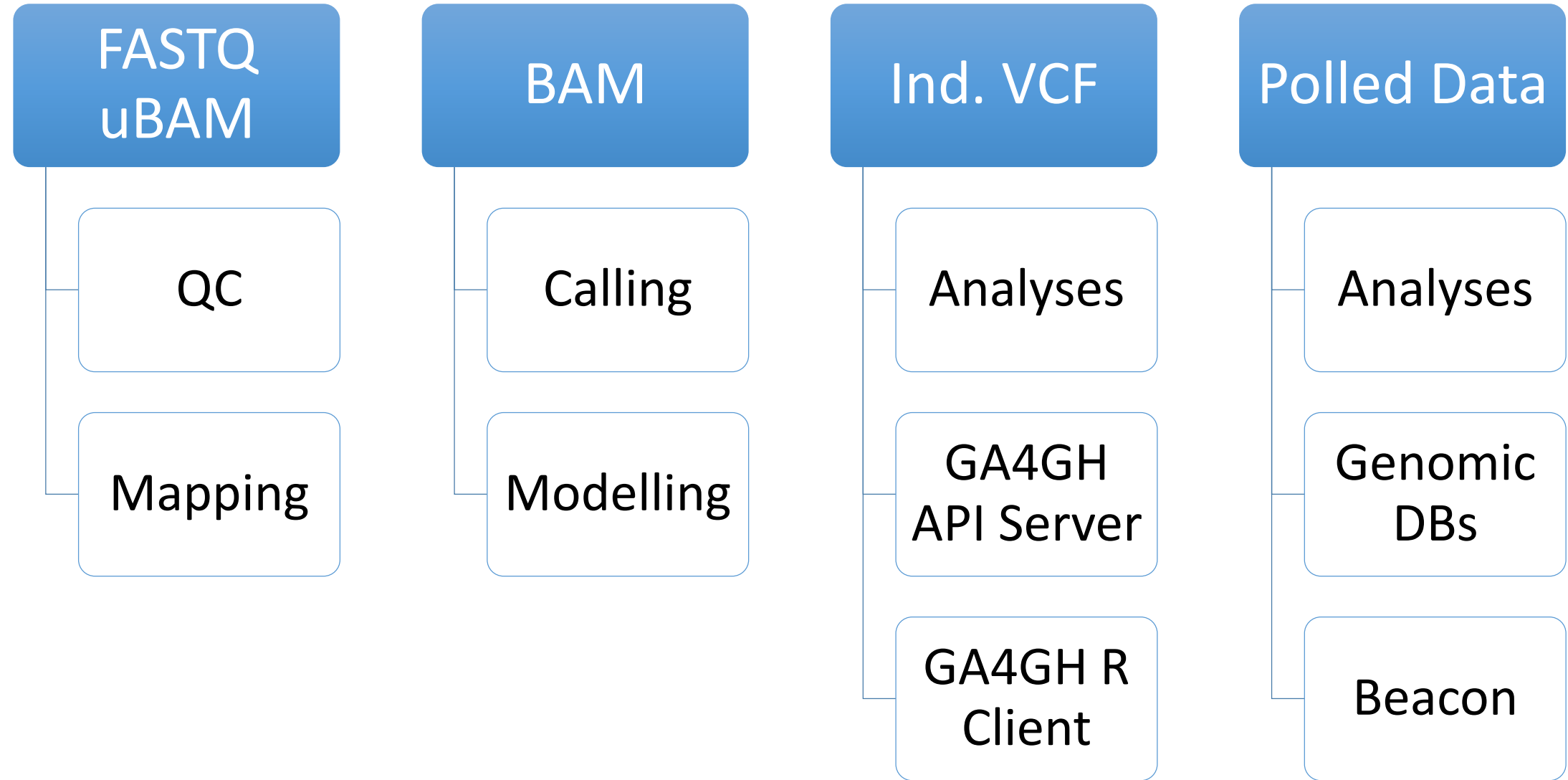| | Chr | Band | Transcripts | Variants | Unique variants | Last |
|---|---|---|---|---|---|---|
| B glycoprotein | 19 | q13.43 | 1 | 0 | 0 | 2015- |
| isense RNA 1 | 19 | q13.43 | 1 | 0 | 0 | 2015- |
| complementation factor | 10 | q21.1 | 10 | 0 | 0 | 2015- |
| nacroglobulin | 12 | p13.31 | 2 | 0 | 0 | 2015- |
| sense RNA 1 (head to head) | 12 | p13.31 | 1 | 0 | 0 | 2015- |
| nacroglobulin-like 1 | 12 | p13 | 4 | 0 | 0 | 2015- |
| nacroglobulin pseudogene 1 | 12 | p13.31 | 1 | 0 | 0 | 2015- |
| -galactosyltransferase 2 | 1 | p35.1 | 1 | 9 | 9 | 2015- |
| -galactosyltransferase | 22 | q13.2 | 10 | 0 | 0 | 2015- |
| -N-acetylglucosaminyltransferase | 3 | p14.3 | 1 | 0 | 0 | 2015- |
| , adrenocortical insufficiency, alacrimia | 12 | q13 | 9 | 0 | 0 | 2015- |
| :yl-CoA synthetase | 12 | q24.31 | 5 | 0 | 0 | 2015- |
| :yl-CoA synthetase pseudogene 1 | 5 | q35 | 1 | 0 | 0 | 2015- |
| mide deacetylase | 3 | q25.1 | 2 | 0 | 0 | 2015- |
| mide deacetylase-like 2 | 3 | q25.1 | 2 | 0 | 0 | 2015- |
| antisense RNA 1 | 3 | q25.1 | 1 | 0 | 0 | 2016- |
| mide deacetylase-like 3 | 1 | p36.21 | 2 | 24 | 24 | 2015- |
| mide deacetylase-like 4 | 1 | p36.21 | 5 | 14 | 14 | 2015- |
| mide deacetylase pseudogene 1 | 3 | q25.1 | 1 | 0 | 0 | 2016- |
| pate aminotransferase | 4 | q33 | 6 | 0 | 0 | 2015- |
| A antioxidant enzyme domain containing 1 | 9 | q22.32 | 3 | 0 | 0 | 2015- |

# Integration with existing tools

- Creating custom-tracks for use with UCSC Genome Browser;
  - All variants;
  - All variants found in at least 1%;
  - All variants found in at least 5%;

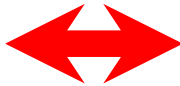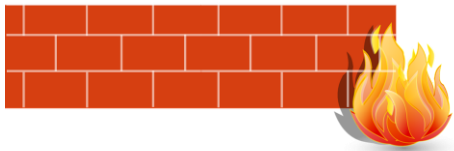- Setting up a local infra-structure to accommodate BIPMed data via UCSC Genome Browser;
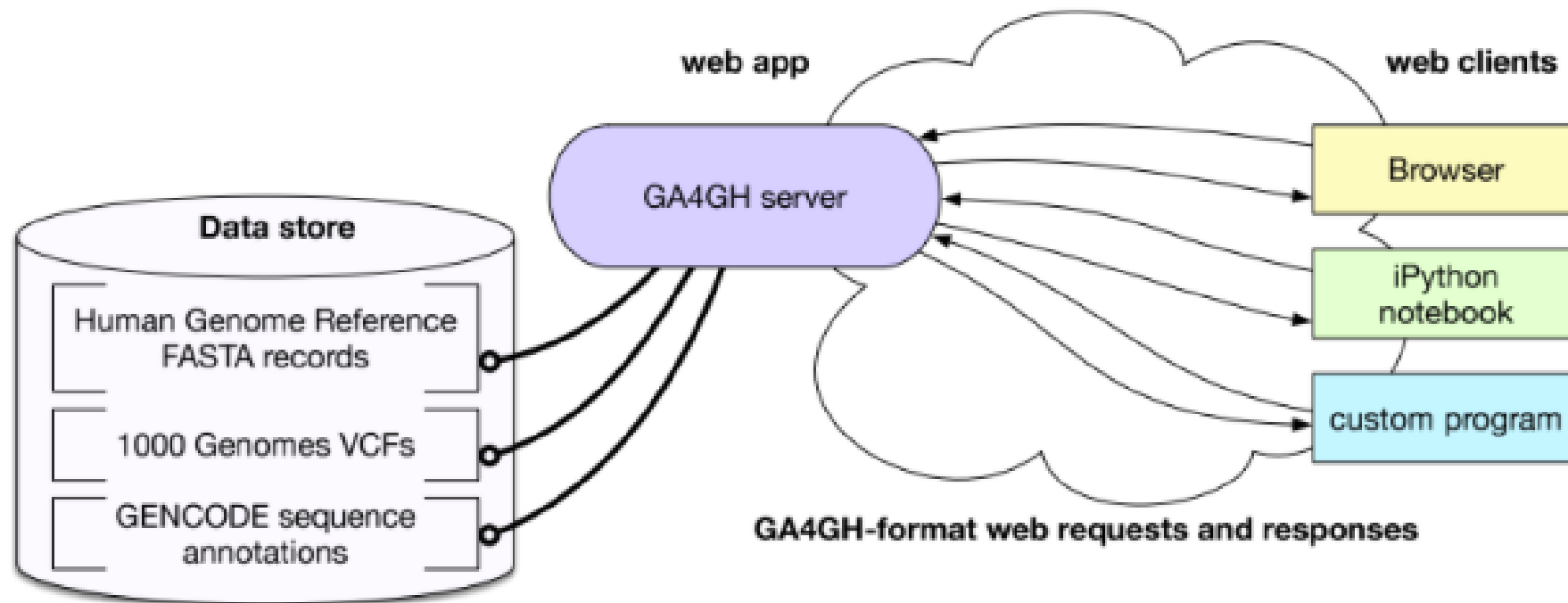
# Data Management – Security – API

# Programmatic Access – Databases

- DB infra-structure uses GA4GH Server API;

- We developed the GA4GH R Client to assist on programmatic access;

- DB access via GA4GH API is possible but not allowed as of today, as authentication is still being developed;

- Tools can be used internally;

# GA4GH Server API

- The ga4gh-server software implements data schema for representing genomic variation data.

- It provides an application program interface (API) through the internet.

- Example: http://1kgenomes.ga4gh.org/

- Project page: https://github.com/ga4gh/server

# GA4GH Server API



Source: GA4GH

# GA4GHclient package

- GA4GHclient is a Bioconductor package for accessing GA4GH API data servers.

- Provides tools to perform programmatic access to genomic data servers that follow the GA4GH standards.

- Integrates with other Bioconductor packages allowing creation of complex genomic data analysis.

- Also provides a web application to interact with genomic data.

- Project page: https://github.com/labbcb/GA4GHclient

GA4GH R Client

# BIPMed Beacon and Beacon Network

# Programmatic Access – Beacon

- BIPMed Beacon connected to Beacon Network;

- Beacon Network can be accessed via Beacon API;

- Beacon API can access BIPMed Beacon directly;

- Tools are under development and designed to allow for queries through computer programs, assisting on variant filtering;

# Programmatic Access – Beacon

```
> getBeaconData(1, 13272, 'C', 'GRCh38', 'bipmed')
[
  {
    "beacon": {
      "id": "bipmed",
      "name": "BIPMed",
      "organization": "Brazilian Initiative on Precision Medicine",
      "description": "Variants identified on Brazilian subjects who belong to the reference population.",
      "aggregator": false,
      "enabled": false,
      "visible": false,
      "createdDate": "2015-11-13",
      "supportedReferences": ["HG38"]
    },
    "query": {
      "chromosome": "CHR1",
      "position": 13272,
      "allele": "C",
      "reference": "HG38"
    },
    "response": true
  }
]
```

# Bottleneck Server

- To avoid genome-wide queries, the BIPMed Beacon uses IP Throttling;

- IP Throttling is performed by the Bottleneck Server, as used by UCSC;

- In practice, the Beacon slows down requests when too many come from the same IP.

# Bottleneck Server

**IP 123.456.789 connects**
- Beacon tells Bottleneck about the user;
- Bottleneck tracks IP;
- Starts counter=0ms

**123.456.789 asks a question;**
- Bottleneck penalizes the client by updating counter to 150ms;
- Client does not ask questions for 1s, counter reduced by 10ms.

**123.456.789 asks more questions;**
- Counter is updated in increments of 150ms;
- If counter gets to 10s, response is delayed in 10s;
- If counter gets to 20s, IP is blocked until counter < 20s.

# BIPMed Website
# http://www.bipmed.org

# Website security

- Website is secured by Cloudfare;
- Contents are copied across many Cloudfare servers around the world;
- Every request to the website is initially sent to the Cloudfare server;
- Cloudfare analyzes requests and decides whether or not they come from a client flagged as a threat;

# Website security

- If client is not a threat:
  - Cloudfare finds the server that is the closest to the client;
  - This server replies to the request by sending cached versions of the website;
  - Connection is much faster and our server does not need to handle traffic;
- If client is a threat:
  - Cloudfare will not reply to requests.
  - Our "original" server is protected from attacks.
- All services provided by us may use a similar system.

# Reliability

- Genomic databases and Beacons are mirrored;
- Current mirrors:
    - USP
    - School of Medicine – UNICAMP;
    - Institute of Chemistry – UNICAMP;
- Mirrors are designed to work as backup systems in case of problems with the main server.

# Thank you!

- Visit our website: http://www.bipmed.org
- Email us: contact@bipmed.org