

# Using Impala For Variant Pipelines

Summer Elasady and Denise Mauldin

11-19-2015

# Connect!

- GitHub tutorials and PowerPoint:

[https://github.com/summerela/impala\\_training/](https://github.com/summerela/impala_training/)

# Agenda

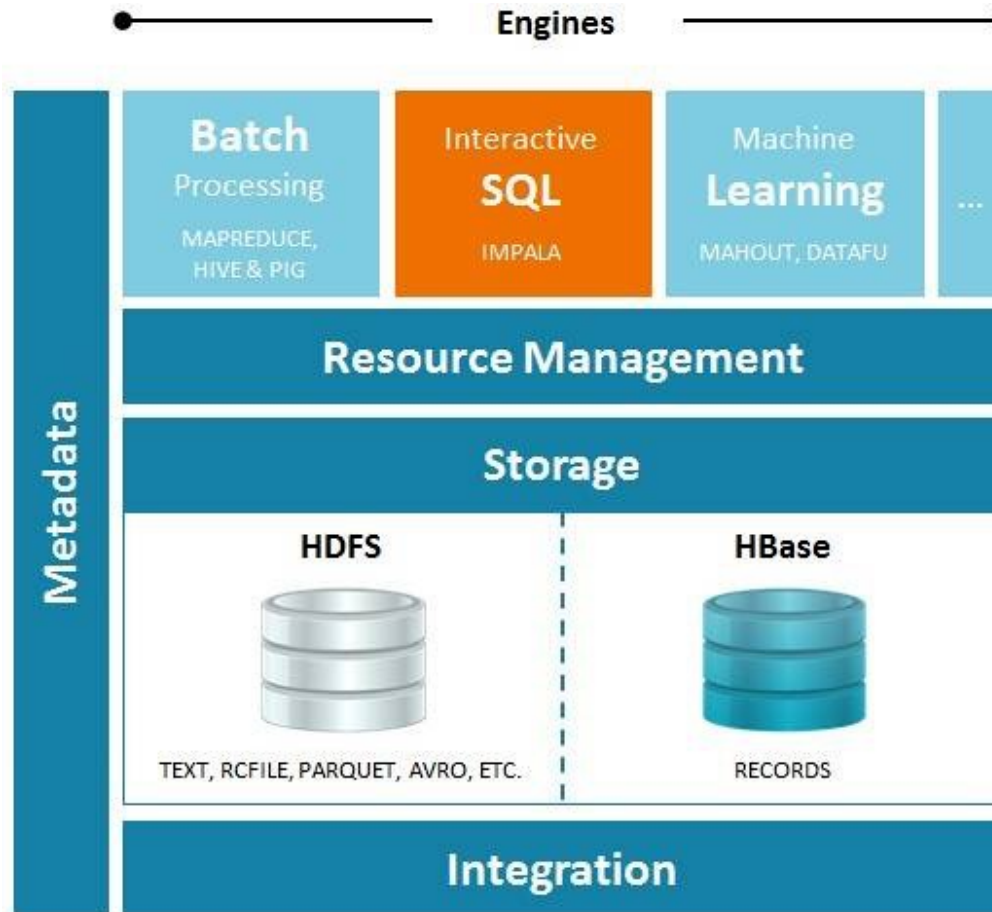
- I. What is Impala?
  - I. Use Cases
  - II. System Structure
  - III. What's available?
- II. Connecting to Impala
  - I. Hue Web Interface
  - II. Impala Shell
  - III. R
  - IV. Python
- III. Creating Queries
- IV. Sample Pipelines
- V. Getting Help



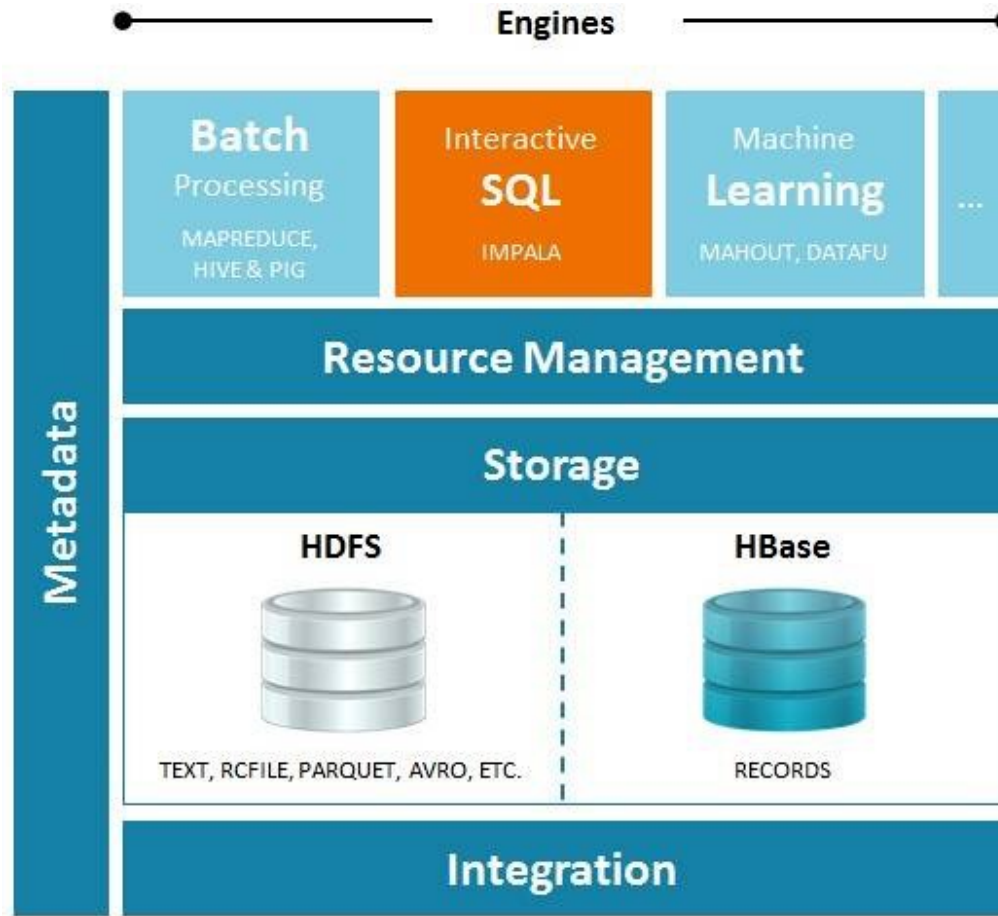
# Use Cases: Impala in the real world

- Locate and annotate variants based on:
  - Pathogenicity
  - Gene region
  - Allele frequency
- Clarity2 Challenge
  - ACMG actionable genes
  - Pathogenic
- Candidate Script
  - Centralized annotation of candidate variants
- Newborn Screening Genes
  - Rare variants
  - Pathogenic

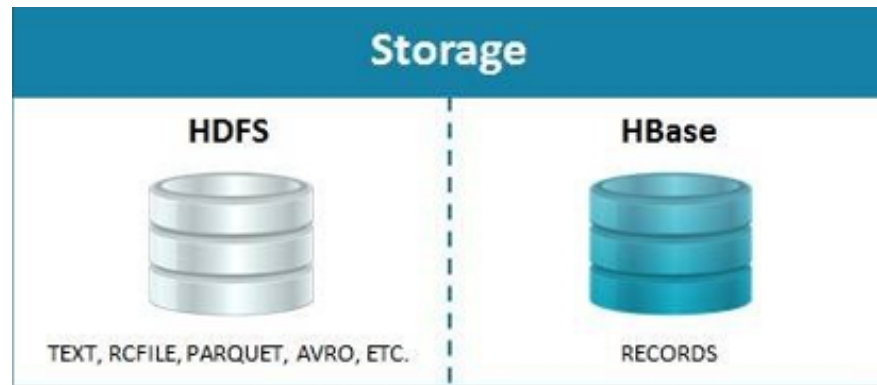
# What exactly is impala?



# What exactly is impala?



# What exactly is impala?

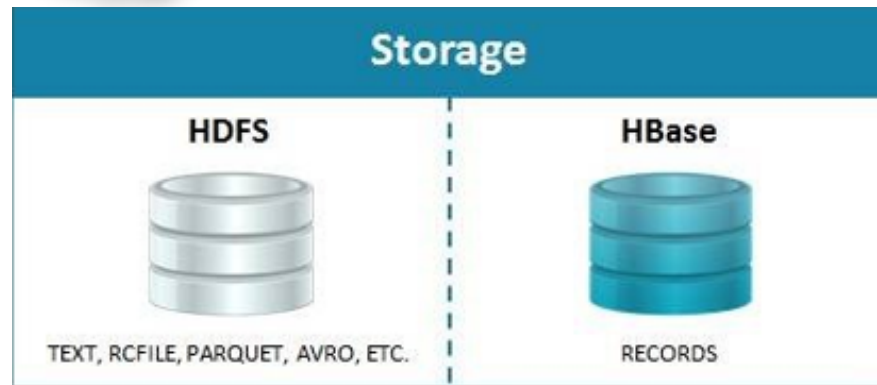


# What exactly is impala?



Store Files:

- Variants
- Reference DB's
- Analysis results





# What exactly is impala?

Raw VCF file transformed and uploaded to HDFS:

```
##fileformat=VCFv4.1
##fileDate=20150819
##source=bin/makeVCF.pl
##reference=file:///proj/famgen/resources/Kaviar-150812-ISB/bin/./tabixedRef/hg19.gz
##version=<Kaviar-150812 (hg19)>
##kaviar_url=<http://db.systemsbiology.org/kaviar>
##publication=<Glusman G, Caballero J, Mauldin DE, Hood L and Roach J (2011) KAVIAR: an accessible system for testing SNV novelty. Bioinformatics, doi: 10.1093/bioinformatics/btr540>
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AC,Number=A,Type=Integer,Description="Allele Count">
##INFO=<ID=AN,Number=1,Type=Integer,Description="Total number of alleles in data sources">
##INFO=<ID=END,Number=.,Type=Integer,Description="End position">
##INFO=<ID=DS,Number=A,Type=String,Description="Data Sources containing allele">
#CHROM POS ID REF ALT QUAL FILTER INFO
1 10001 . T C . . AF=0.0000384;AC=1;AN=26028;DS=SS6004475
1 10002 . A C,T . . AF=0.0001153,0.0000384;AC=3,1;AN=26028;DS=HGDP00927|HGDP00998|HGDP01284,HGDP01029
1 10002 . A AT . . AF=0.0000384;AC=1;AN=26028;DS=HGDP00521
1 10003 . A C,T . . AF=0.0000384,0.0000768;AC=1,2;AN=26028;DS=HGDP01284,HGDP00521|HGDP00927
1 10004 . C A . . AF=0.0000384;AC=1;AN=26028;DS=HGDP01284
1 10018 . C T . . AF=0.0000384;AC=1;AN=26028;DS=HGDP00998
1 10019 rs775809821 TA T . AF=0.0000384;AC=1;AN=26028;END=10020;DS=MaIay
```

# What exactly is impala?

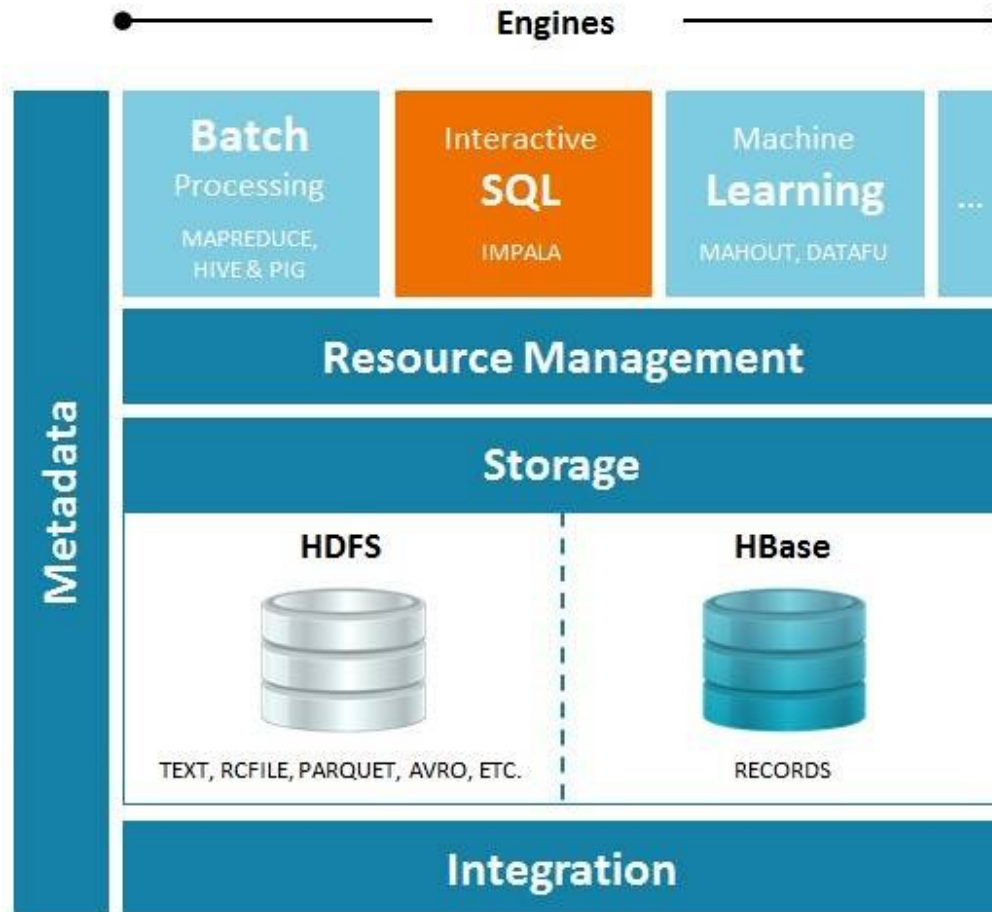
Read as a table by impala:

Data sample for kaviar

[View in Metastore Browser](#)

chrom	pos	stop	rs_id	ref	alt	qual	filter	allele_freq	allele_cnt	allele_num	sources
10	98917931	None	rs75578462	C	T	None	None	0.053711399436	1398	26028	ADNI GMIK9 GS000009930 GS0000117
10	98917981	None	rs545362632	C	T	None	None	0.000230499994359	6	26028	phase3-MSL
10	98917985	None	rs17112469	T	C	None	None	0.106577500701	2774	26028	!Gub ADNI Desmond Tutu GMIK6 GMIA
10	98918018	None	rs746271737	G	A	None	None	3.84000013582e-05	1	26028	UK10K
10	98918073	None	rs756450057	T	C	None	None	3.84000013582e-05	1	26028	UK10K

# What exactly is impala?



# What exactly is impala?



# What exactly is impala?



## SQL – Structured query language

SELECT columns  
FROM database.table\_name  
WHERE parameters

# What exactly is impala?

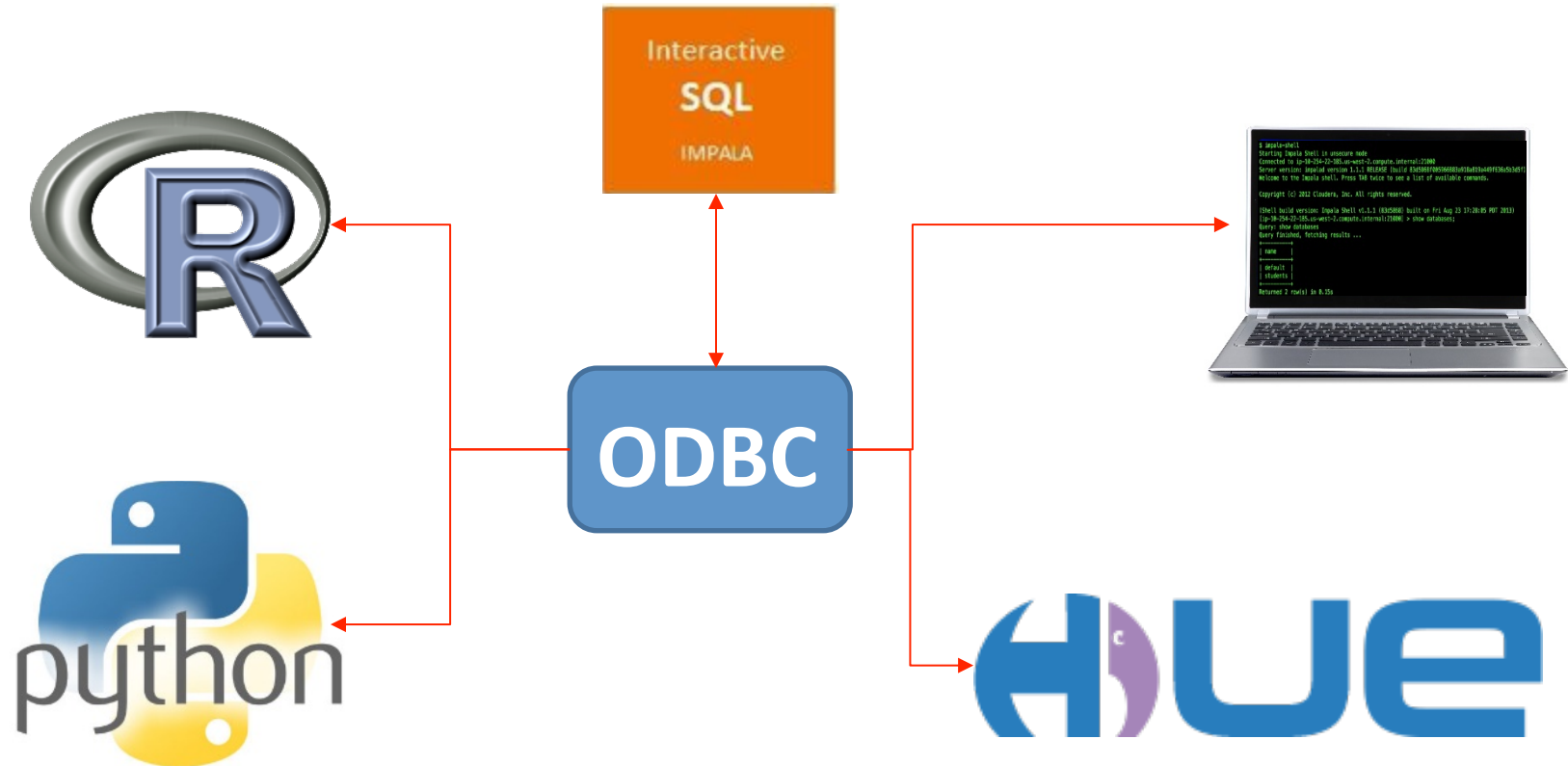


Interact with impala to:

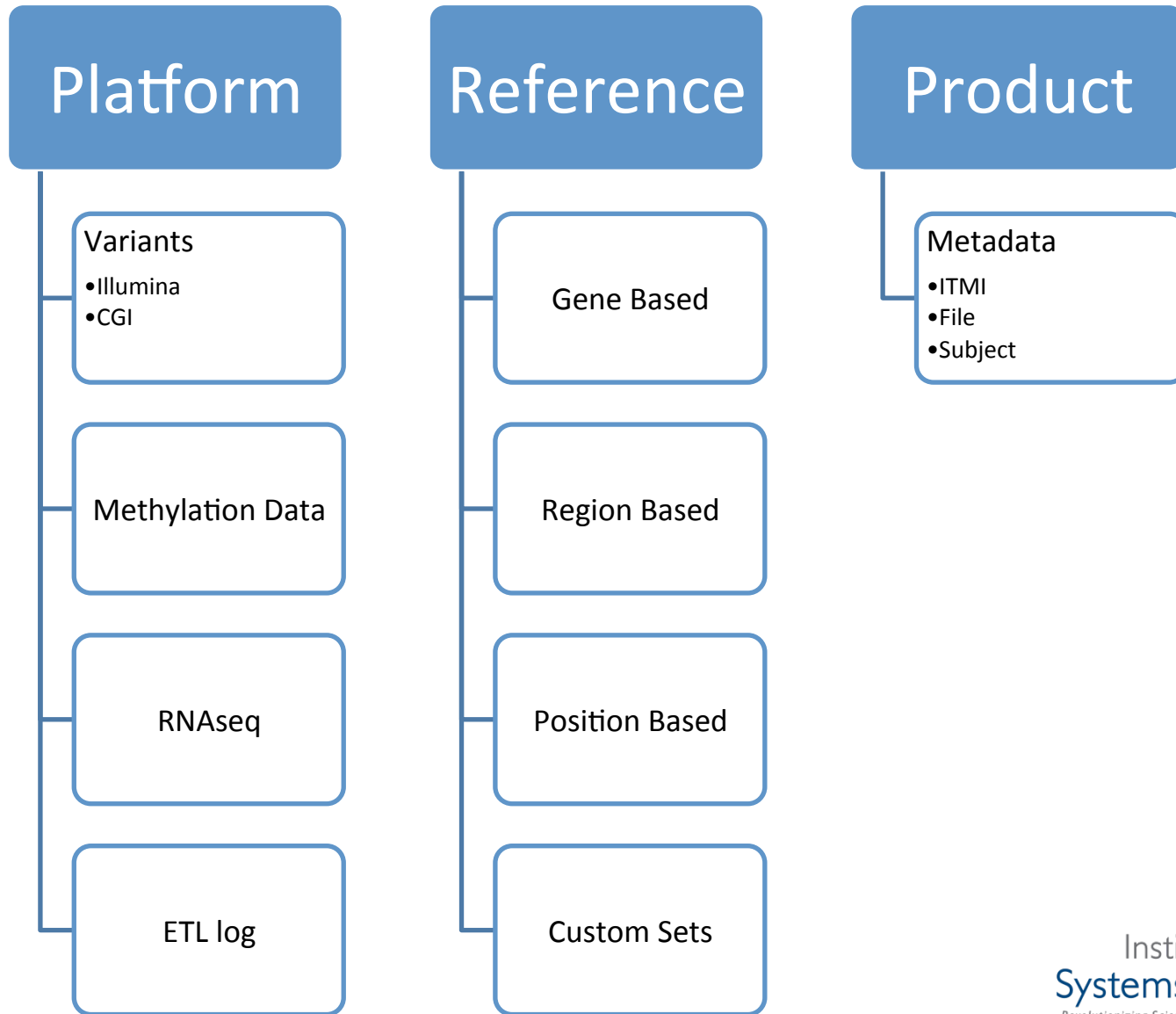
- Filter
- Locate
- Annotate



# What exactly is impala?



# Database Structure





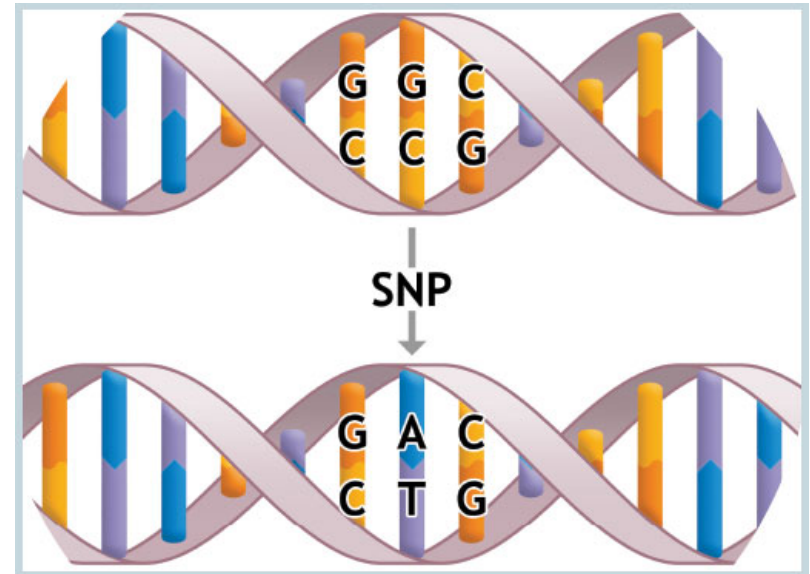
# Benefits

- Consistency
  - Normalized Data
  - 1-based coordinates
  - Universal Column Names
- Speed
  - Quickly run queries across thousands of genomes

# What's Available: Variants

Complete  
genomics 

illumina®

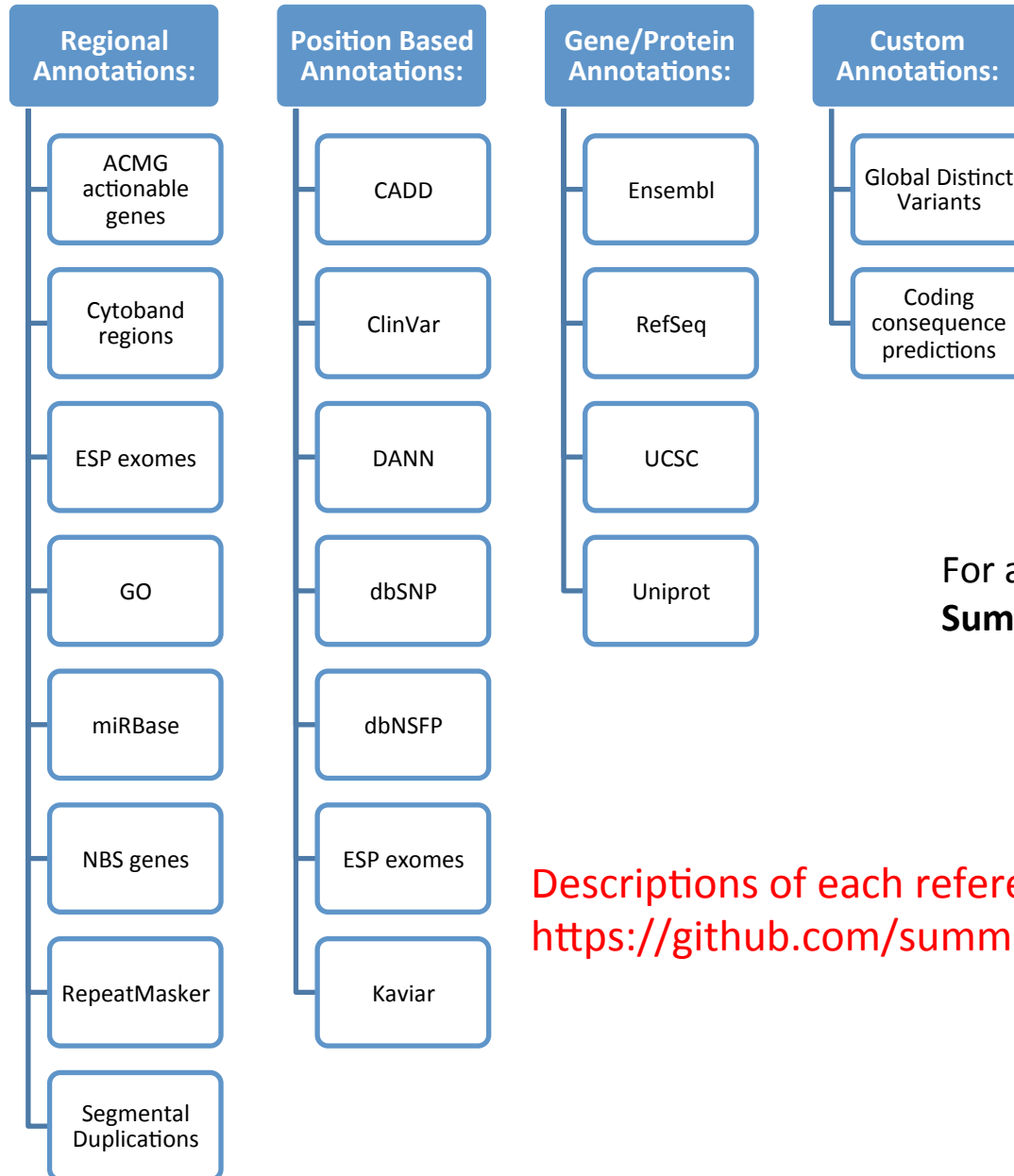


101

102

103

# Reference Sources



For additional annotations email:  
**Summer.Elasady@systemsbiology.org**

Descriptions of each reference:  
<https://github.com/summerela/external-reference-data-catalog>

# Accessing Impala



# Accessing Impala



# Accessing Impala

**cloudera**<sup>®</sup>  
IMPALA



# Accessing Impala

cloudera<sup>®</sup>  
IMPALA



## Pros:

- Easy, web based, quick
- Easy to do query testing and EXPLAIN plan

## Cons:

- Cuts off results inconsistently at 100,000 rows
- AJAX connection issues:
  - Connection drops
  - Canceling queries

HUE

# Accessing Impala





# Accessing Impala





# Accessing Impala



# Accessing Impala



# Accessing Impala

**cloudera**<sup>®</sup>  
IMPALA



ODBC  
driver

# Accessing Impala

**cloudera**<sup>®</sup>  
IMPALA



ODBC  
driver

R ODBC



Impyla  
Ibis



# Hue Interface



## ISB Impala Login:

- Login to your impala web interface

## Hue Tutorial:

[https://github.com/summerela/  
impala\\_training/blob/master/using\\_hue.pdf](https://github.com/summerela/impala_training/blob/master/using_hue.pdf)

# Connecting with impala-shell

- Tutorial:

[https://github.com/summerela/impala\\_training/blob/master/impala\\_shell.ipynb](https://github.com/summerela/impala_training/blob/master/impala_shell.ipynb)



# Connecting with Python

- Python:

[https://github.com/summerela/impala\\_training](https://github.com/summerela/impala_training)

– Click on the ‘launch binder’ icon

The icon consists of two adjacent rectangular buttons. The left button is dark grey with the word 'launch' in white lowercase text. The right button is pink with the word 'binder' in white lowercase text.

– Click on connect\_python.ipynb

# Connecting with R

- R:
  - Launch R:
    - Launch your R web interface
  - Tutorial:  
[https://github.com/summerela/impala\\_training/blob/master/connect\\_with\\_R.md](https://github.com/summerela/impala_training/blob/master/connect_with_R.md)
  - Script:  
[https://github.com/summerela/impala\\_training/blob/master/connect\\_R.R](https://github.com/summerela/impala_training/blob/master/connect_R.R)

# Building Queries

- Choose your interface
- Follow along (and copy/paste queries):  
[https://github.com/summerela/impala\\_training/blob/master/building\\_queries.md](https://github.com/summerela/impala_training/blob/master/building_queries.md)

# Basic Pipeline

1. Connect to impala
2. Input genes of interest
3. Input subject id's of interest
4. Annotate
5. Filter
6. Export results

[https://github.com/summerela/impala\\_training/blob/master/variant\\_pipeline\\_python.ipynb](https://github.com/summerela/impala_training/blob/master/variant_pipeline_python.ipynb)