

ENCODE DCC data modeling and metadata standardization vignettes

Idan Gabdank
ENCODE Data Coordination Center

GSC20

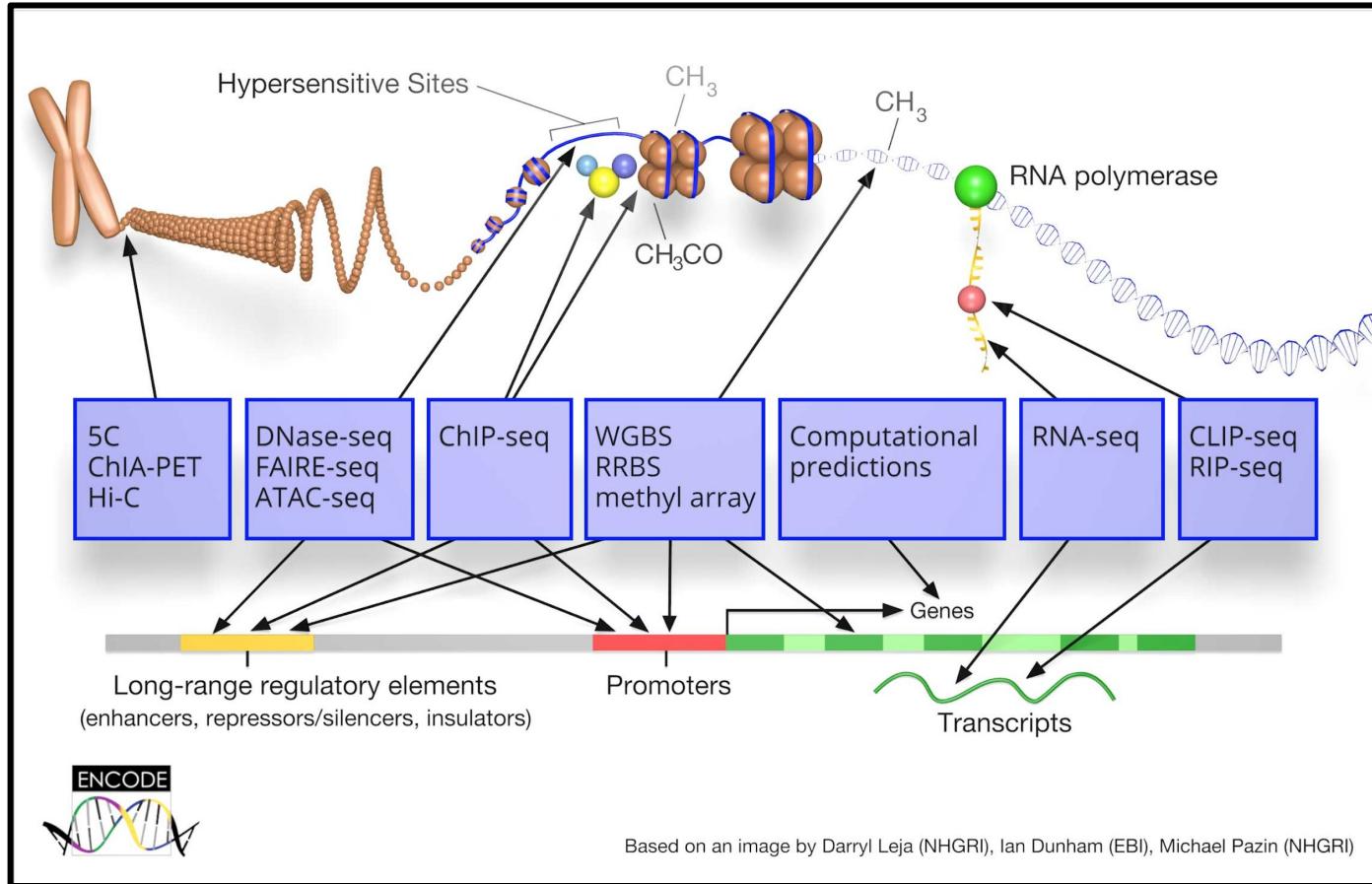
A long journey to reproducible results

Replicating our work took four years and 100,000 worms but brought surprising discoveries, explain **Gordon J. Lithgow, Monica Driscoll and Patrick Phillips.**

How do you turn irreproducible
experiment into reproducible?

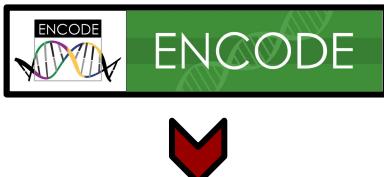
- * methods standardization
- * experimental replication
- * documentation

Encyclopedia of DNA Elements



Encyclopedia of DNA Elements

2003-2007



ENCODE



2007-2011



ENCODE
PHASE 2



2012-2016



ENCODE
PHASE 3

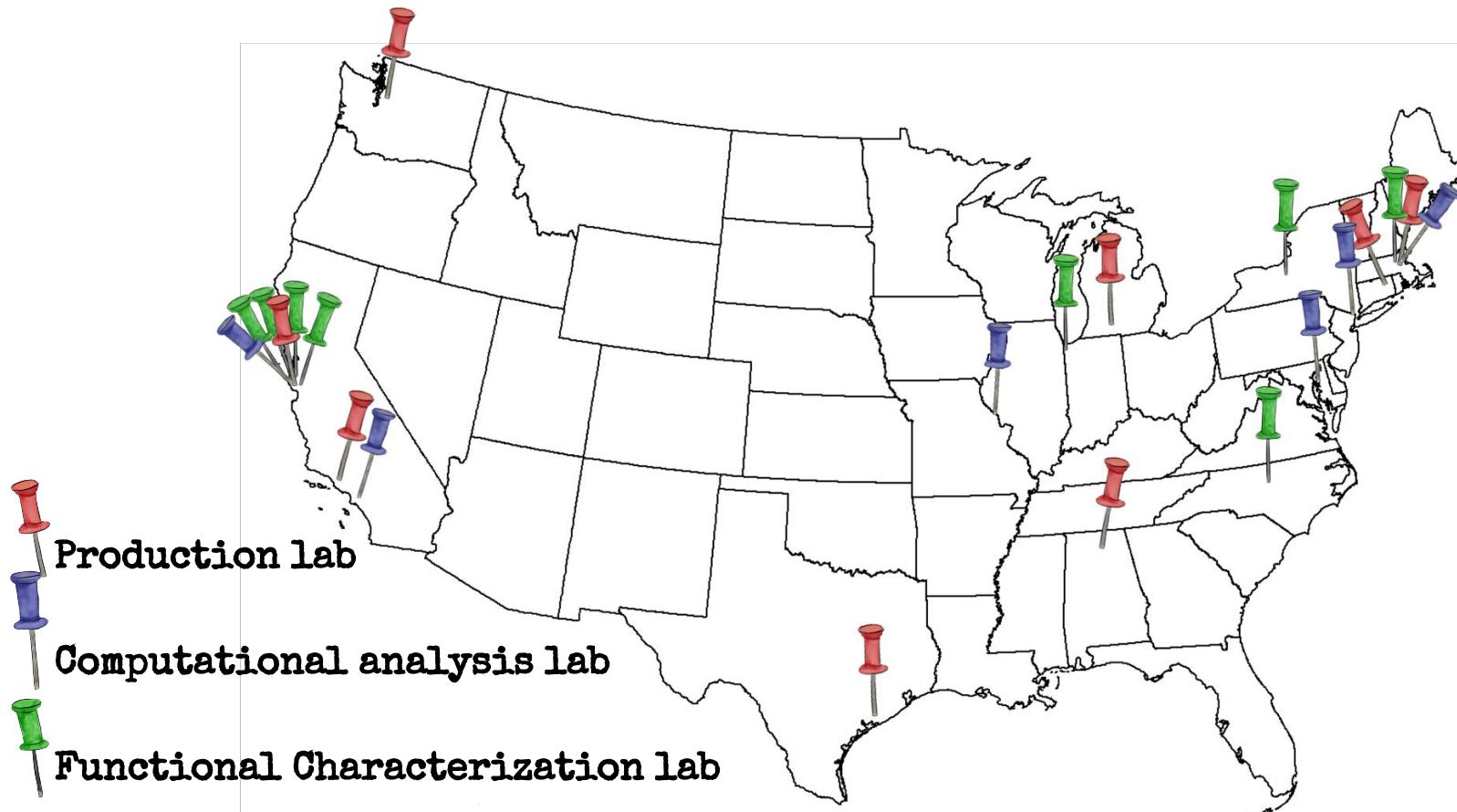
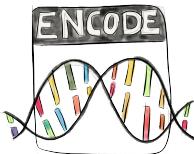


2017-2021

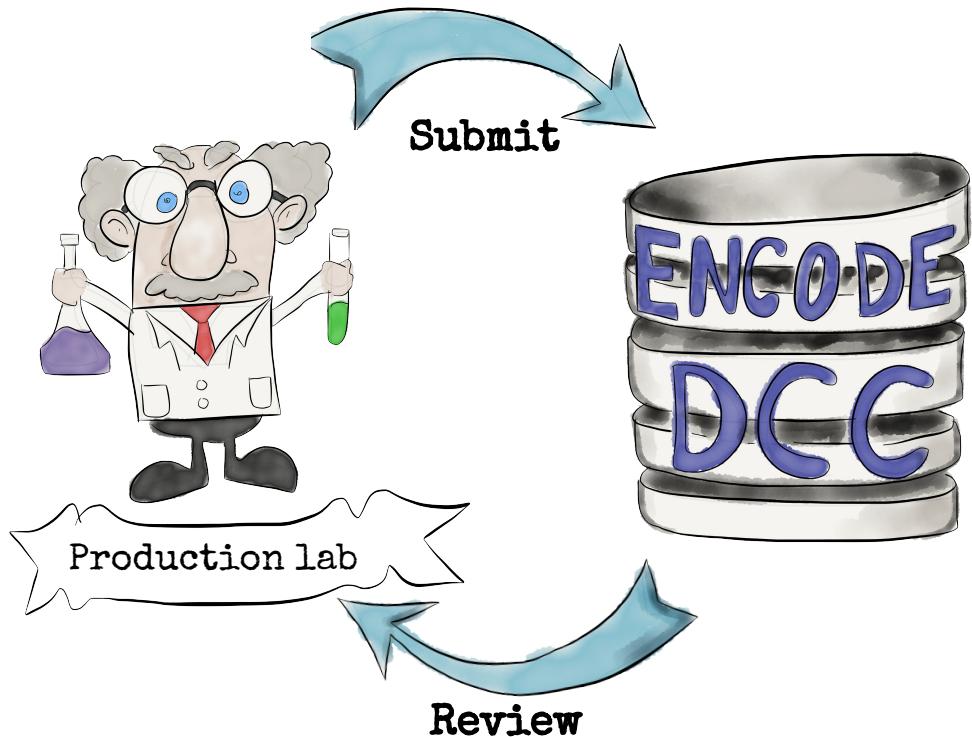


ENCODE
PHASE 4

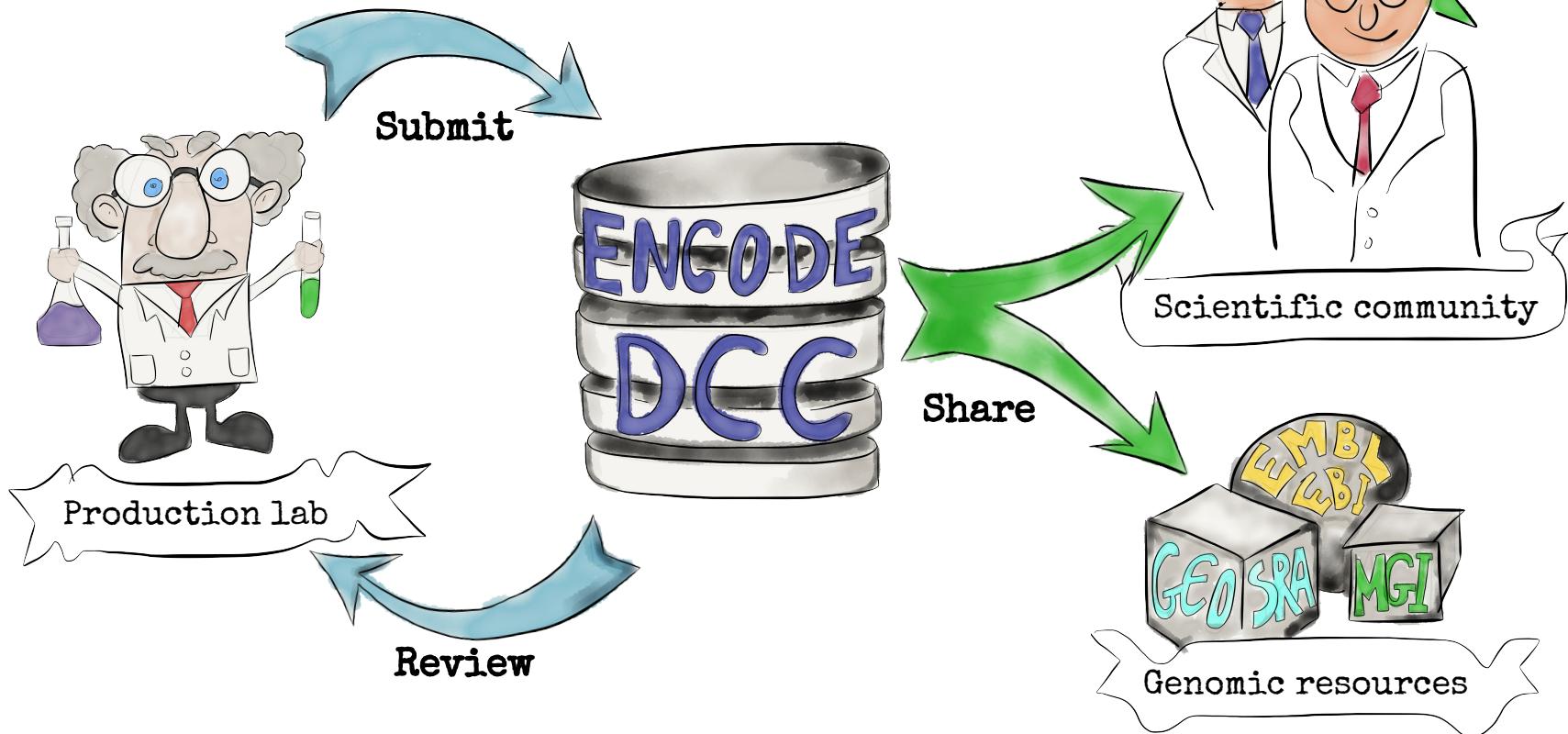
Encyclopedia of DNA Elements (phase IV)

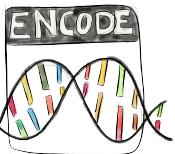


ENCODE DCC role

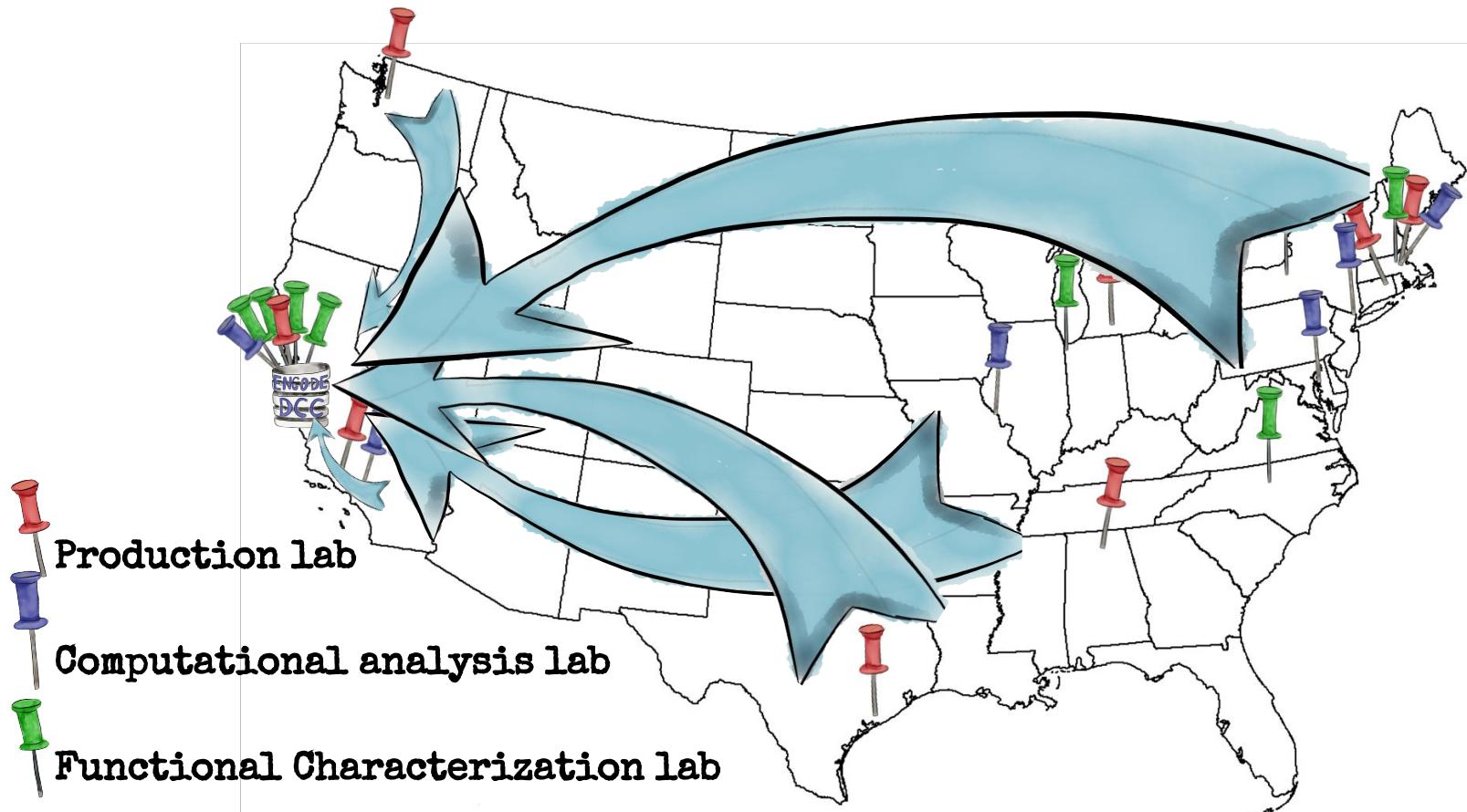


ENCODE DCC role

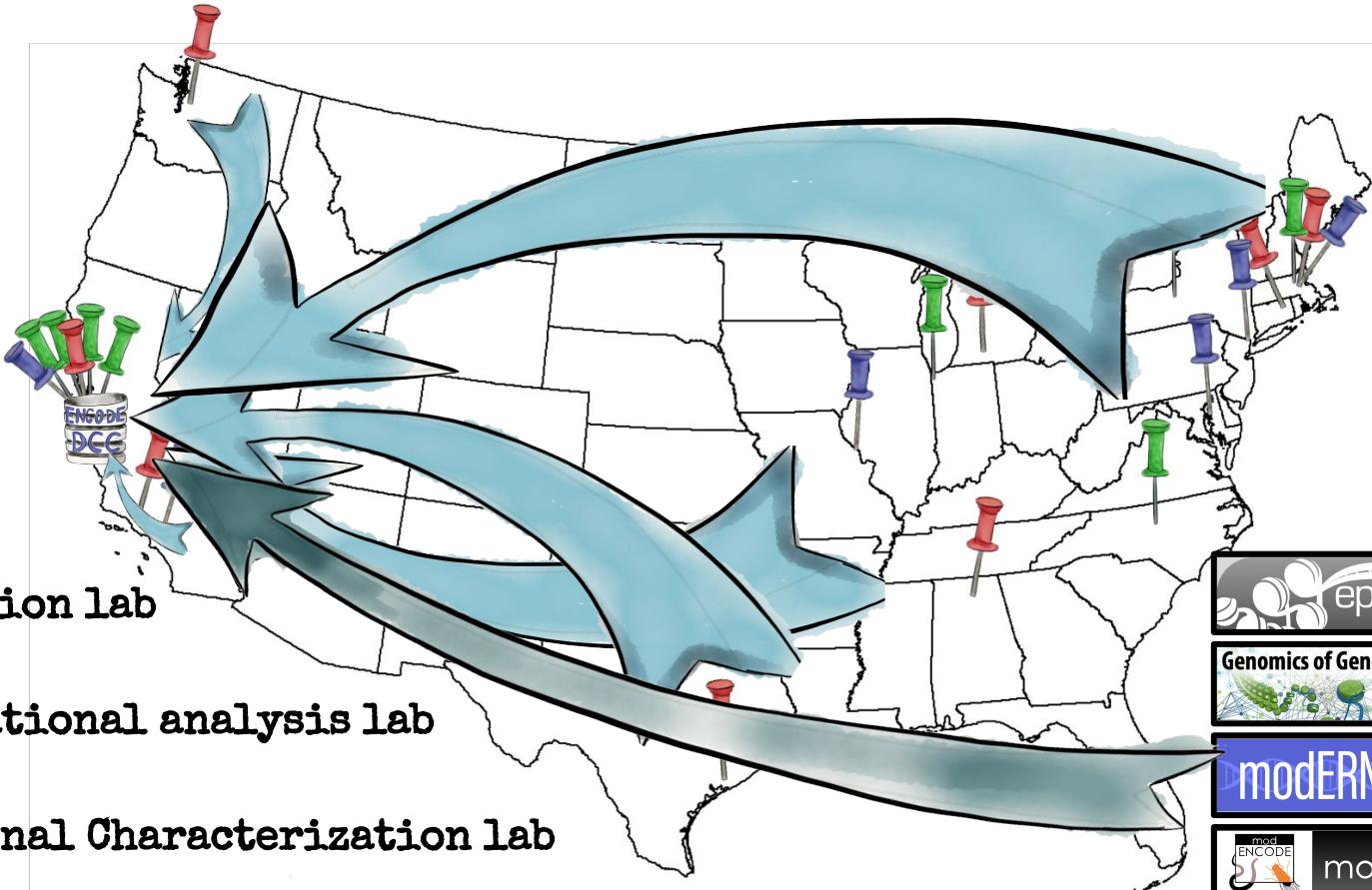
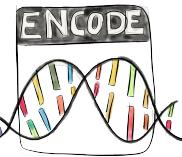




Data flow



Data flow

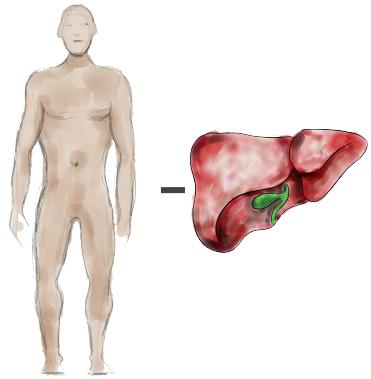


From the lab to the portal



Donor

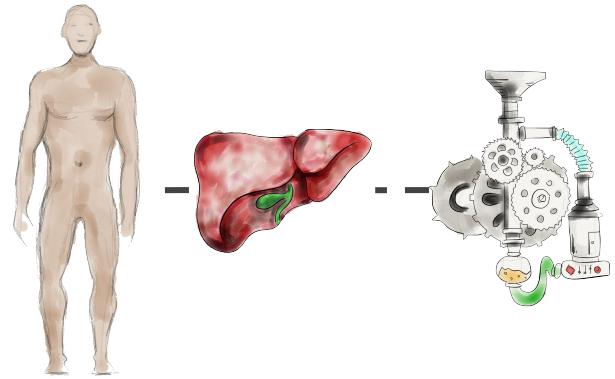
From the lab to the portal



Donor

Sample

From the lab to the portal

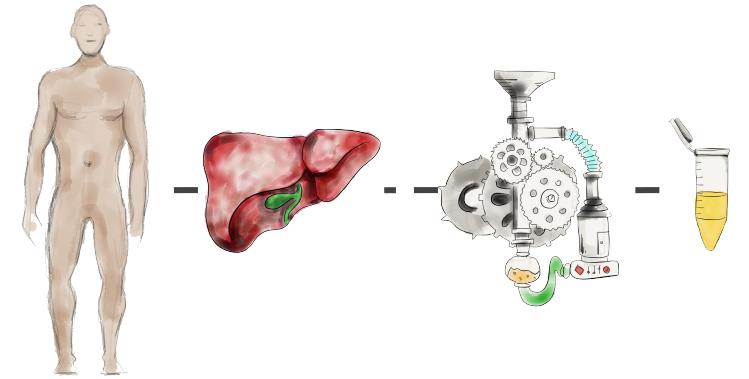


Donor

Sample

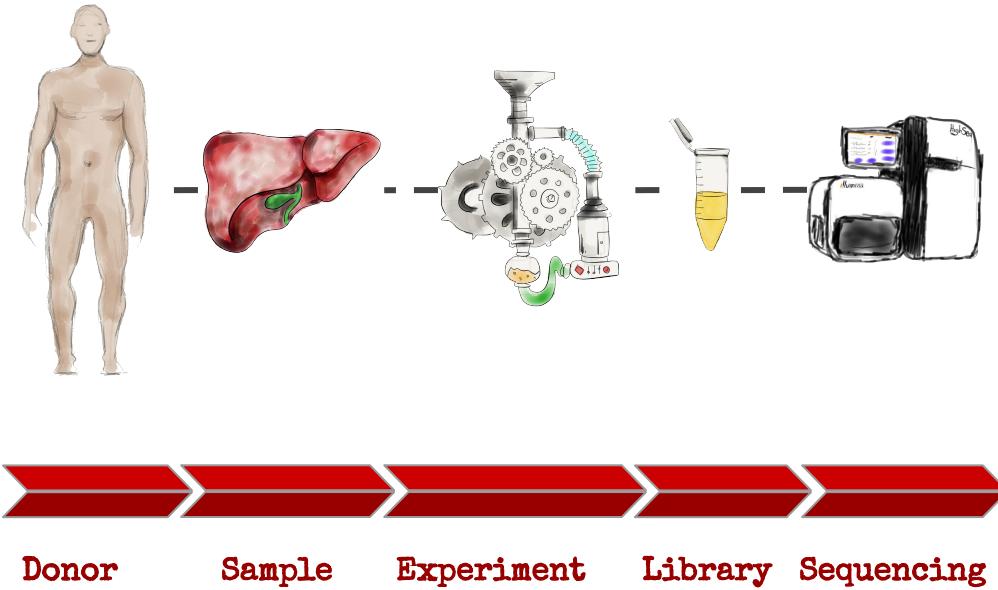
Experiment

From the lab to the portal

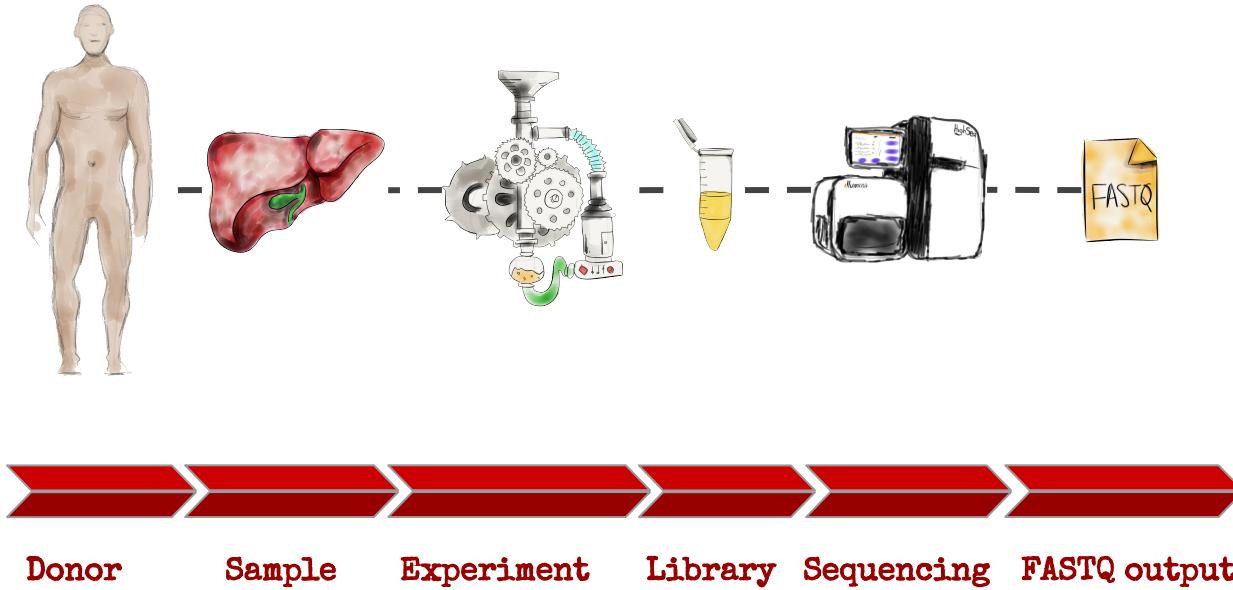


Donor Sample Experiment Library

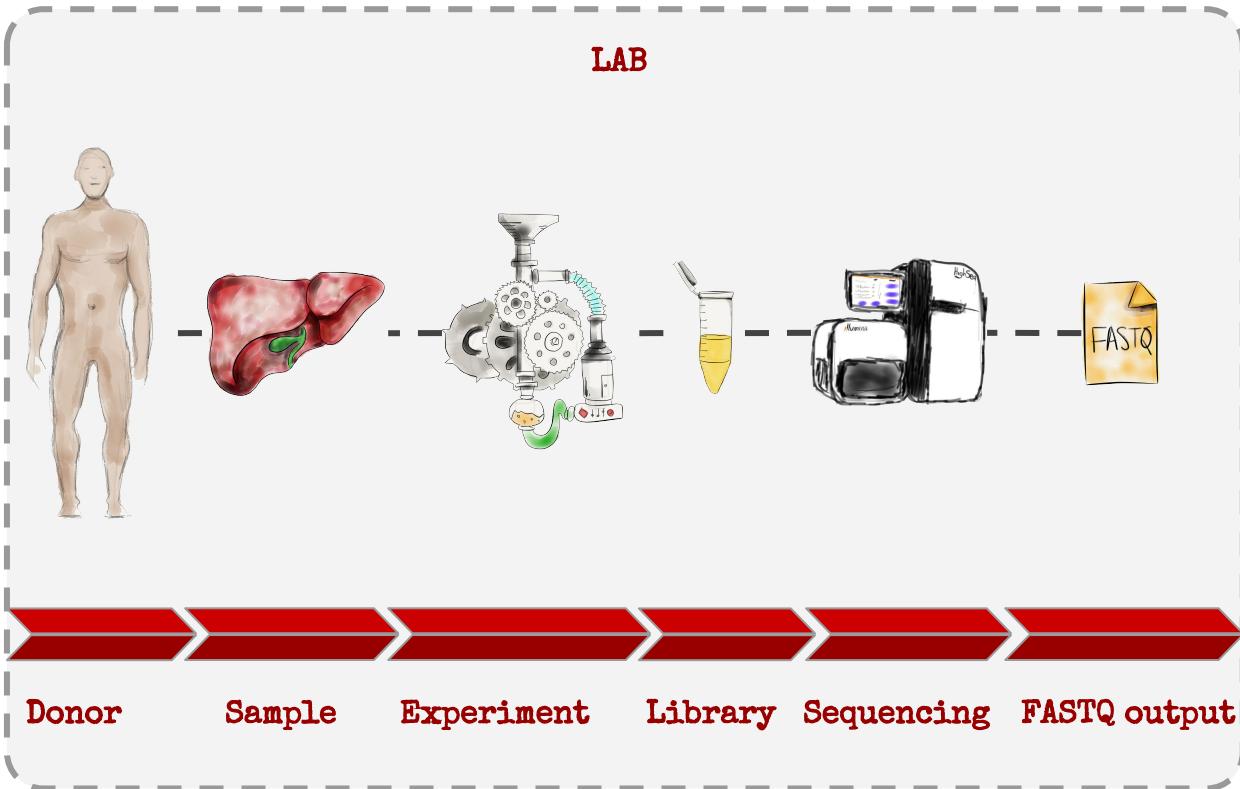
From the lab to the portal



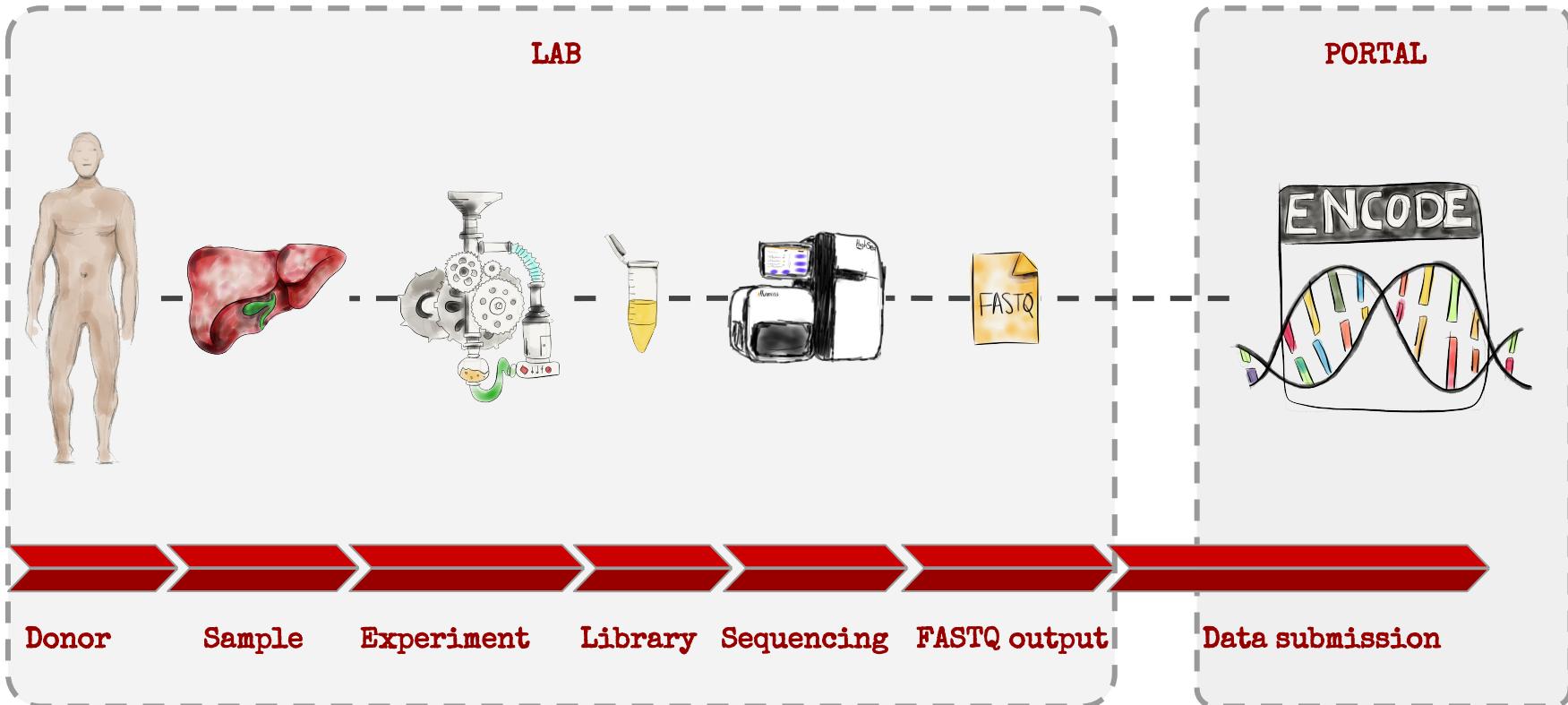
From the lab to the portal



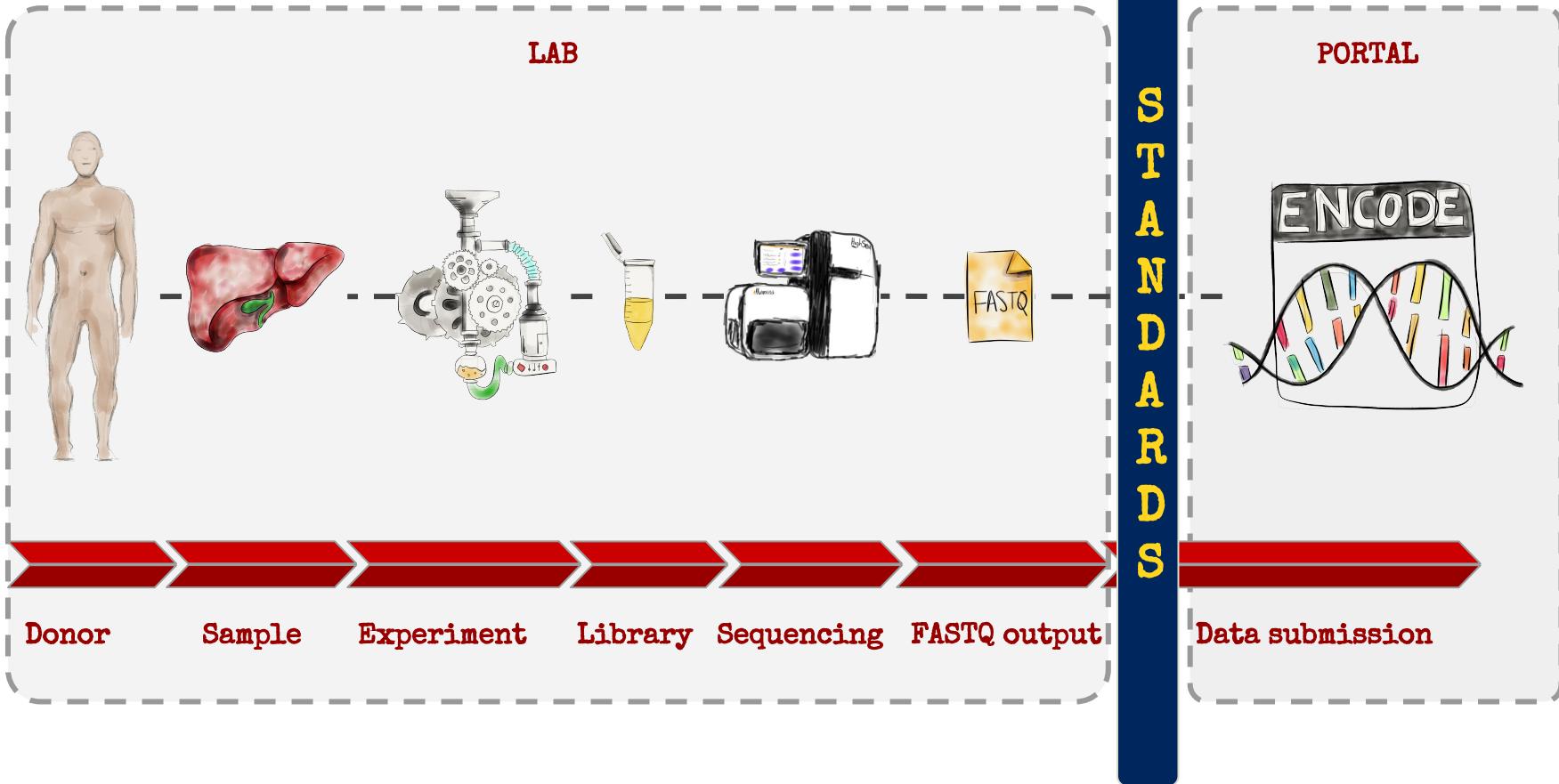
From the lab to the portal



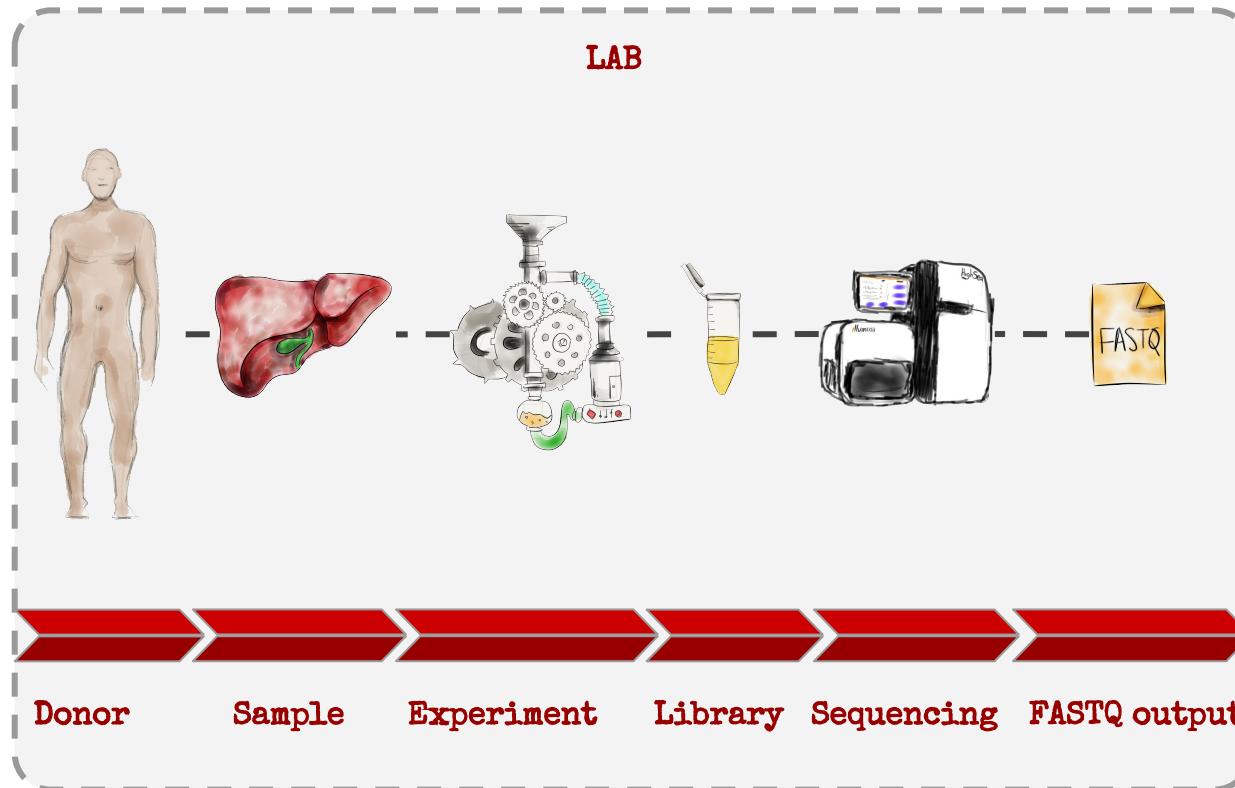
From the lab to the portal



From the lab to the portal



Data categories (JSON objects)



Donor

- * female
- * 51 years old



HumanDonor.json

ENCDO271OUW

Status:
released

Donor information

Accession: ENCDO271OUW

Aliases: gtex:ENC-DEJ, gtex:PT-1LVAN, gtex:ENC-004, bradley-bernstein:Donor GTEX-1LVAN

Donor external identifiers: [GEO:SAMN05897787](#)

Species: *Homo sapiens*

Life stage: Adult

Age: 51 year

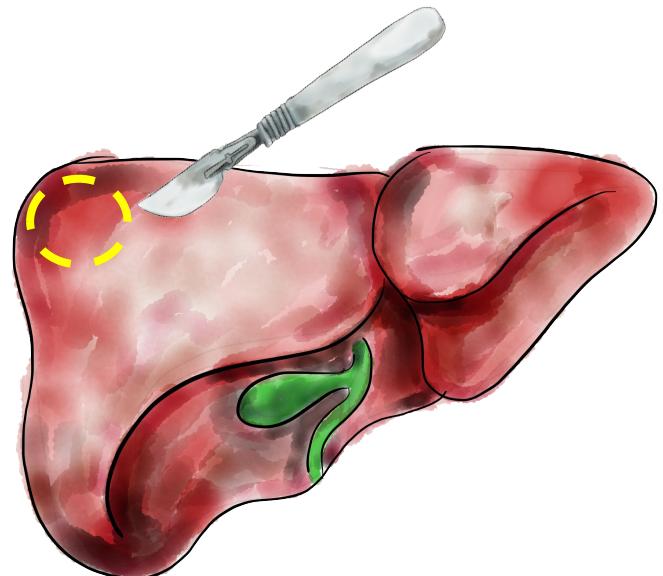
Sex: Female

Tags:



Tissue sample

* right lobe of liver



Biosample.json

ENCBS150EDI / tissue

Status: released

Summary

Term name: right lobe of liver

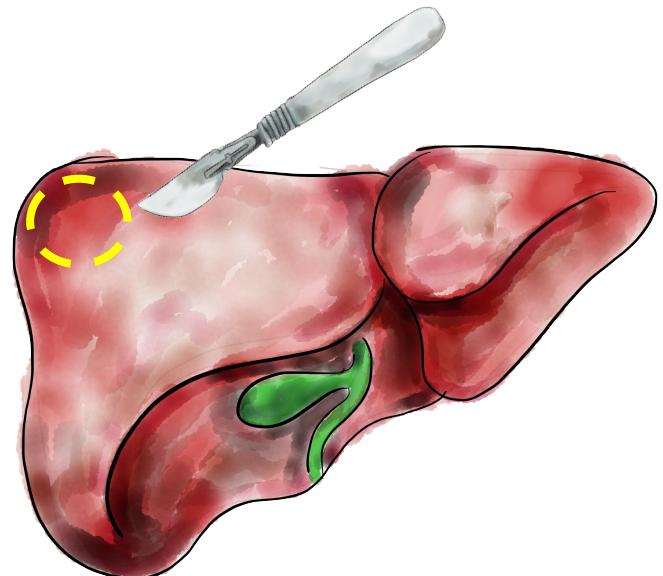
Term ID: UBERON:0001114 [↗](#)

Summary: *Homo sapiens* female adult (53 years) right lobe of liver tissue

Parent of biosamples:

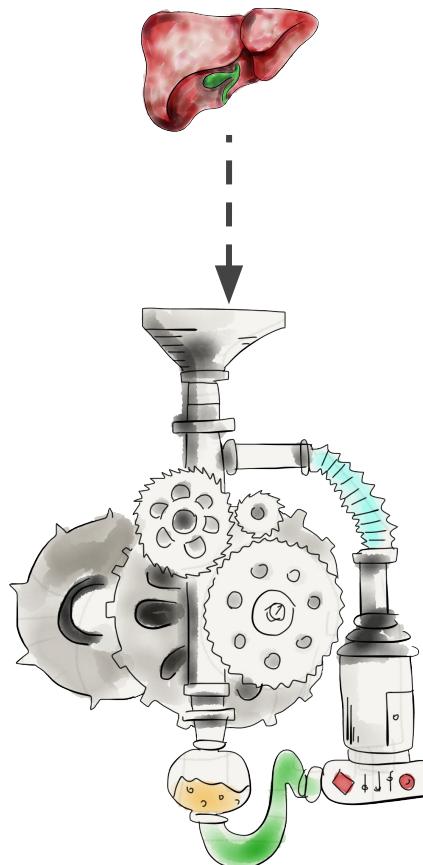
ENCBS220FFU, ENCBS688NPW, ENCBS594VBH, ENCBS796JPW, ENCBS773RMK, ENCBS021IAC, ENCBS904HZU, ENCBS290WWG, ENCBS984AKV, ENCBS682CWZ, ENCBS536THV, ENCBS405OIZ, ENCBS054RZZ, ENCBS034NFF

UBERON:0001114

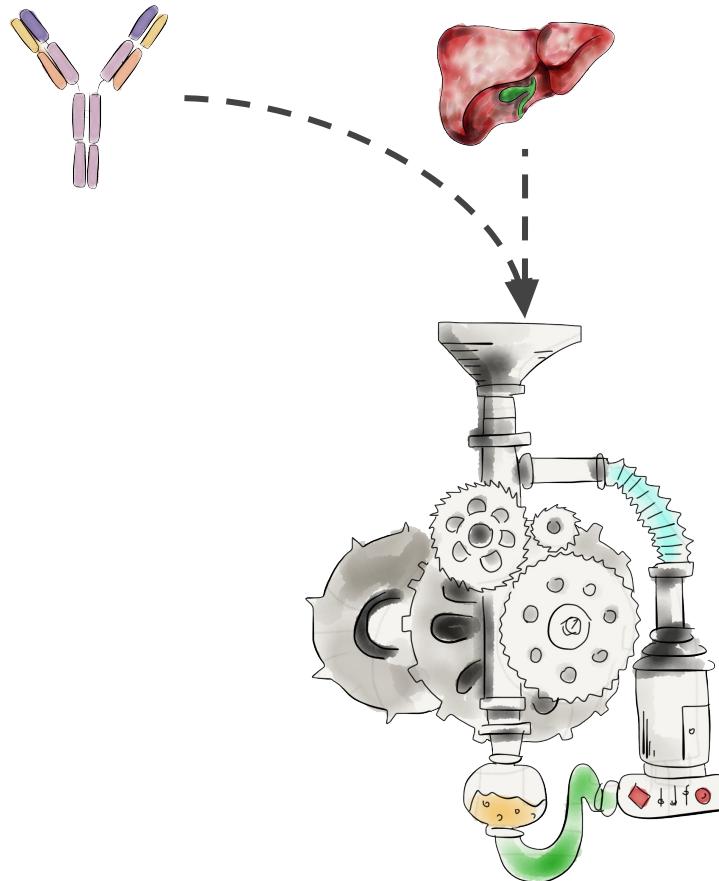


Experiment

CTCF ChIP-seq
OBI:0000716

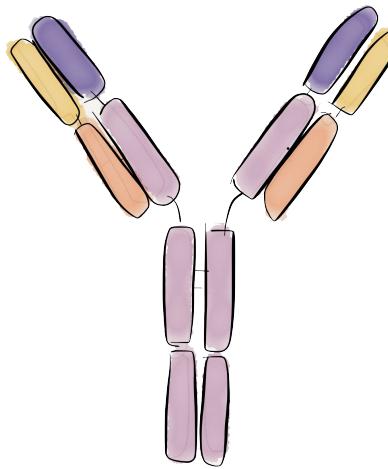


Experiment

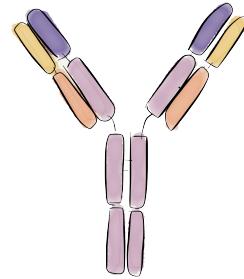


Antibody

- * Vendor
- * Product ID
- * Lot ID



AntibodyLot.json



ENCAB635SXP

Antibody against *Homo sapiens* CTCF

Status: released

2 1

Homo sapiens at least one cell type or tissue

Awaiting
characterization

Source (vendor): ABclonal ↗

Product ID: A1133 ↗

Lot ID: A1133

Targets: CTCF (*Homo sapiens*)

Host: Rabbit

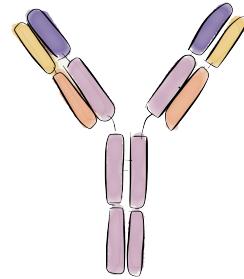
Clonality: Polyclonal

Purification: Affinity

Isotype: IgG

Antigen description: Recombinant protein of human CTCF

AntibodyLot.json



ENCAB635SXP

Antibody against *Homo sapiens* CTCF

Status: released 2 1

Homo sapiens at least one cell type or tissue Awaiting characterization

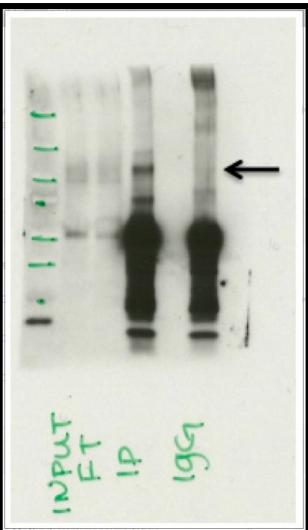
Source (vendor):	ABclonal ↗
Product ID:	A1133 ↗
Lot ID:	A1133
Targets:	CTCF (<i>Homo sapiens</i>)
Host:	Rabbit
Clonality:	Polyclonal
Purification:	Affinity
Isotype:	IgG
Antigen description:	Recombinant protein of human CTCF

Antibody standards

Primary characterization

Secondary characterization

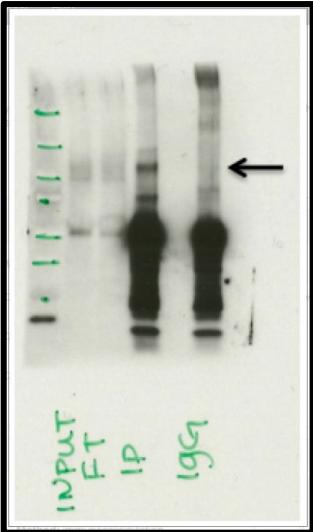
Antibody standards

Primary characterization	Secondary characterization
Western blot or IP Western  A Western blot image showing protein bands across four lanes. The lanes are labeled from left to right: INPUT, FT, IP, and IgG. The IP lane shows a prominent dark band at the same position as the IgG lane, indicated by a black arrow pointing to the right. The INPUT and FT lanes show faint bands. The IgG lane shows a very strong, dark band. INPUT FT IP IgG	

Antibody standards

Primary characterization

Western blot or IP Western

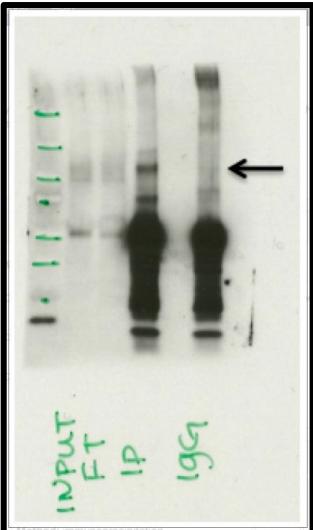
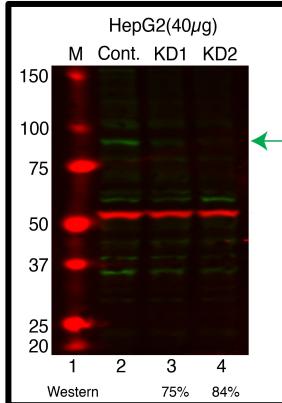


Secondary characterization

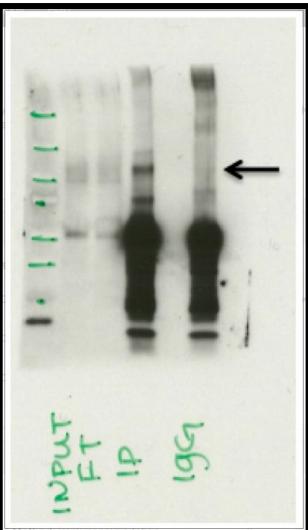
- Mass spectrometry

Spectrum	Name of Protein	Count of Peptides	Ratio(CTCF/IgG Control)
CTCF Band A	Heterogeneous nuclear ribonucleoprotein M (Fragment)	15	NOT IN CONTROL IP
CTCF Band A	C4A protein	13	NOT IN CONTROL IP
CTCF Band A	Isoform 2 of Nucleolar RNA helicase 2	10	NOT IN CONTROL IP
CTCF Band A	Prelinin-1/C	10	NOT IN CONTROL IP
CTCF Band A	Probable ATP-dependent RNA helicase DDX5	10	NOT IN CONTROL IP
CTCF Band A	Transcriptional repressor CTCF	10	NOT IN CONTROL IP
CTCF Band A	Heat shock cognate 71 kDa protein	9	NOT IN CONTROL IP
CTCF Band A	Isoform 2 of Probable ATP-dependent RNA helicase DDX17	9	NOT IN CONTROL IP
CTCF Band A	N-acetyltransferase 10	8	NOT IN CONTROL IP
CTCF Band A	Splicing factor 3B subunit 3	8	NOT IN CONTROL IP
CTCF Band A	SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily A (Fragment)	8	NOT IN CONTROL IP
CTCF Band A	Isoform short of Heterogeneous nuclear ribonucleoprotein U	7	NOT IN CONTROL IP
CTCF Band A	X-ray repair complementing defective repair in Chinese hamster cells 6 (Ku autoantigen, Fragment)	7	NOT IN CONTROL IP
CTCF Band A	Nucleolar protein 5G (Fragment)	6	NOT IN CONTROL IP
CTCF Band A	ATP-dependent RNA helicase A	5	NOT IN CONTROL IP
CTCF Band A	Interleukin enhancer-binding factor 3	5	NOT IN CONTROL IP
CTCF Band A	Isoform UBFL2 of Nucleolar transcription factor 1	5	NOT IN CONTROL IP
CTCF Band A	X-ray repair cross-complementing protein 5	5	NOT IN CONTROL IP
CTCF Band A	DNA topoisomerase 1	4	NOT IN CONTROL IP

Antibody standards

Primary characterization	Secondary characterization
<p>Western blot or IP Western</p> 	<ul style="list-style-type: none">• Mass spectrometry• siRNA or shRNA against the mRNA of the target protein 

Antibody standards

Primary characterization	Secondary characterization
Western blot or IP Western 	<ul style="list-style-type: none">• Mass spectrometry• siRNA or shRNA against the mRNA of the target protein• ChIP-seq data from a previously characterized antibody•••

Antibody standards

Primary characterization ✓

Secondary characterization ✓



characterized to standards

Homo sapiens

C4-2B, HepG2, K562, RWPE2, LNCAP, GM12878, HCT116, 22Rv1, VCaP, RWPE1, HeLa-S3,
Panc1, MCF-7

Antibody standards

Primary characterization ✓

Secondary characterization ✓

CHALLENGE

● characterized to standards

Homo sapiens

C4-2B, HepG2, K562, RWPE2, LNCAP, GM12878, HCT116, 22Rv1, VCaP, RWPE1, HeLa-S3,
Panc1, MCF-7

Experiment.json

Experiment summary for ENCSR911GFJ

Status: released 2

Summary		Attribution	
Assay:	ChIP-seq	Lab:	Michael Snyder, Stanford
Target:	CTCF	Award:	U54HG006996 (Michael Snyder, Stanford)
Biosample summary:	<i>Homo sapiens</i> right lobe of liver female adult (53 years)	Project:	ENCODE
Biosample Type:	tissue	External resources:	GEO:GSE105829
Replication type:	unreplicated	Aliases:	michael-snyder:ChIPss-674
Description:	CTCF ChIP-seq on human right lobe of liver	Date submitted:	December 6, 2016
Nucleic acid type:	DNA	Date released:	December 13, 2016
Size range:	450-650	Tags:	
Strand specificity:	Non-strand-specific		
Platform:	Illumina HiSeq 4000		
Controls:	ENCSR336OPU		



Replicates

Biological replicate	Technical replicate	Summary	Biosample	Antibody	Library
1	1	female adult (53 years) right lobe of liver tissue	ENCBS904HZU	ENCAB830JLB	ENCLB331LXY

Raw data file(s)

- * file format and integrity
- * data duplication prevention

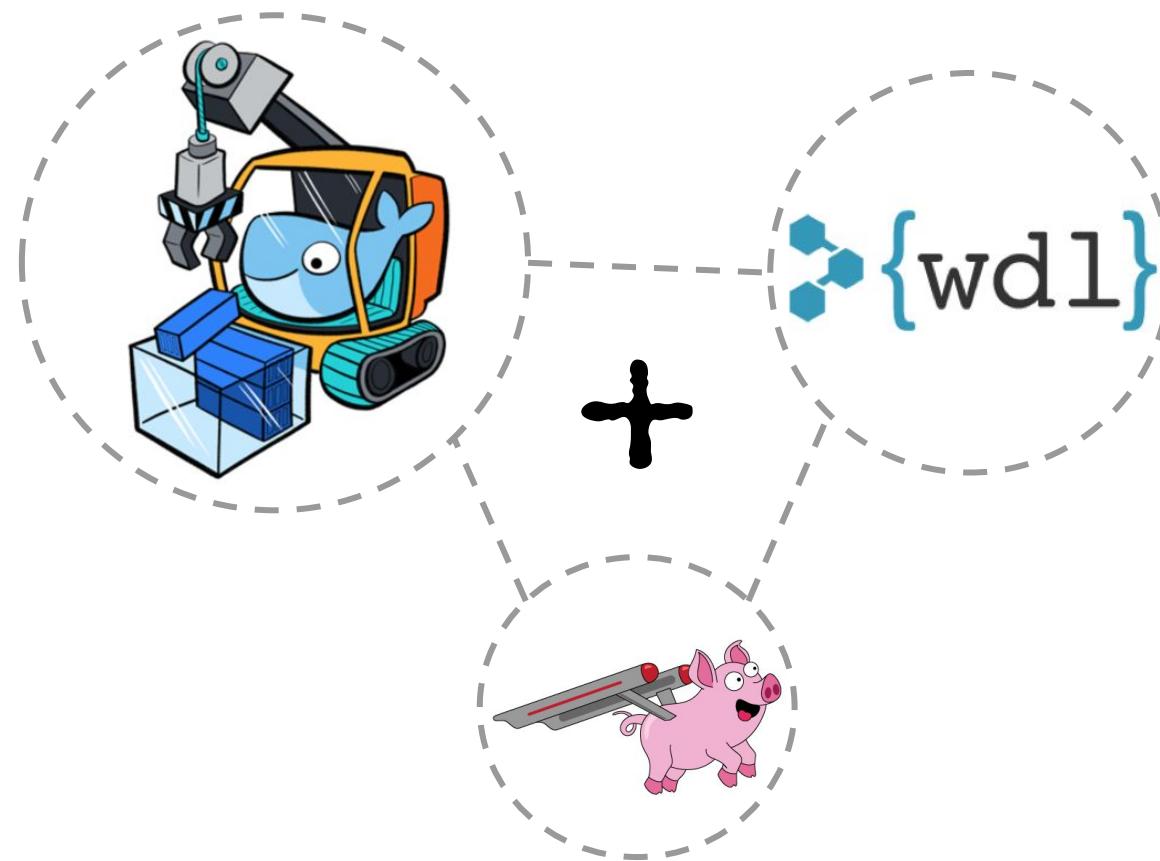


Raw data file(s)

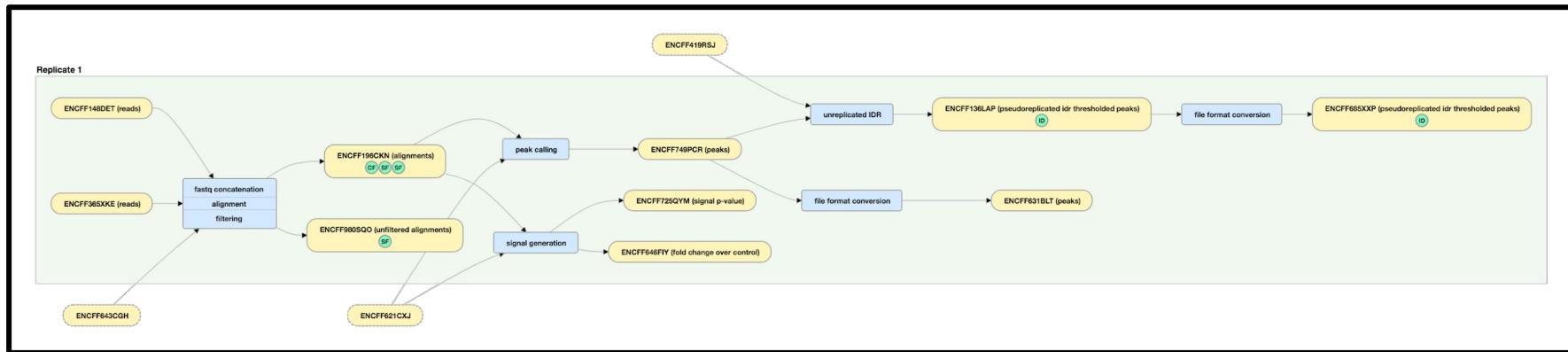
- * file format and integrity
- * data duplication prevention



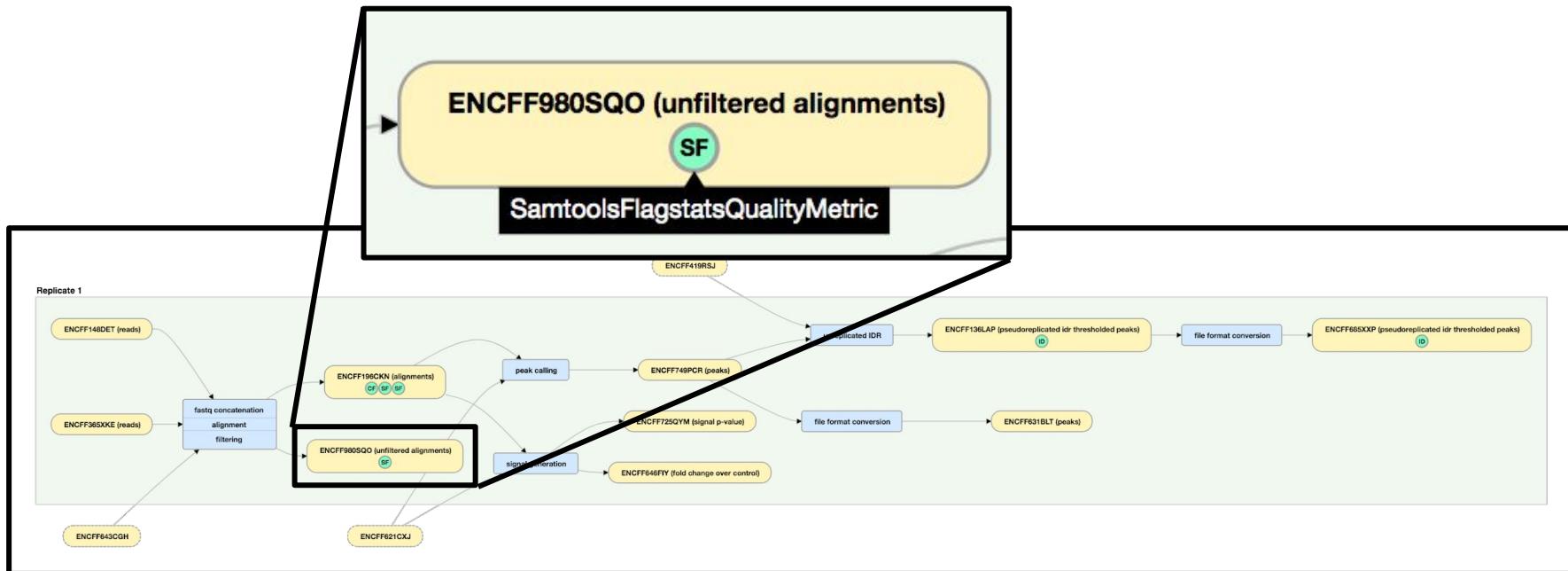
Uniform processing pipelines



Analysis file graph



QC metrics



Audits

- * metadata integrity
- * ENCODE consortium data standards

Badges

RED for ERROR 

ORANGE for NOT_COMPLIANT 

YELLOW for WARNING 

GREY for INTERNAL DCC 

Badges

Experiment summary for ENCSR000DKV

Status: released

Internal: unreviewed

1 1 2

	Insufficient read length	
	Inconsistent platforms	
	Missing derived_from	
	Biological replicates with identical biosample	

Badges

Experiment summary for ENCSR000DKV

Status: released

Internal: unreviewed

1 1 2



Insufficient read length

Fastq file [/files/ENCFF000ROU/](#) has read length of 35bp. For mapping accuracy ENCODE standards recommend that sequencing reads should be at least 50bp long. (See [/data-standards/chip-seq/](#))

Fastq file [/files/ENCFF000ROX/](#) has read length of 33bp. For mapping accuracy ENCODE standards recommend that sequencing reads should be at least 50bp long. (See [/data-standards/chip-seq/](#))

Fastq file [/files/ENCFF000ROZ/](#) has read length of 35bp. For mapping accuracy ENCODE standards recommend that sequencing reads should be at least 50bp long. (See [/data-standards/chip-seq/](#))



Inconsistent platforms



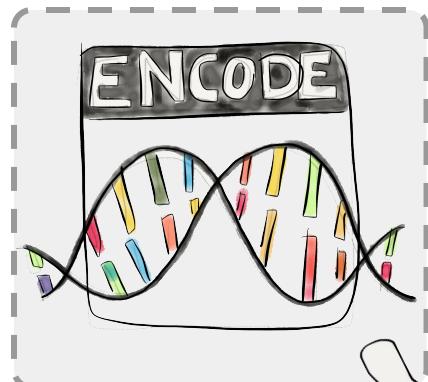
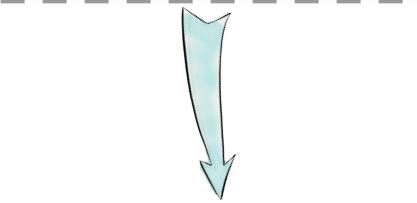
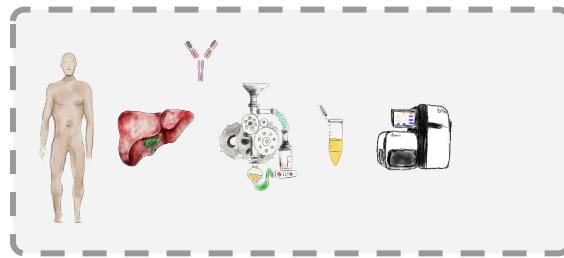
Missing derived_from



Biological replicates with identical biosample

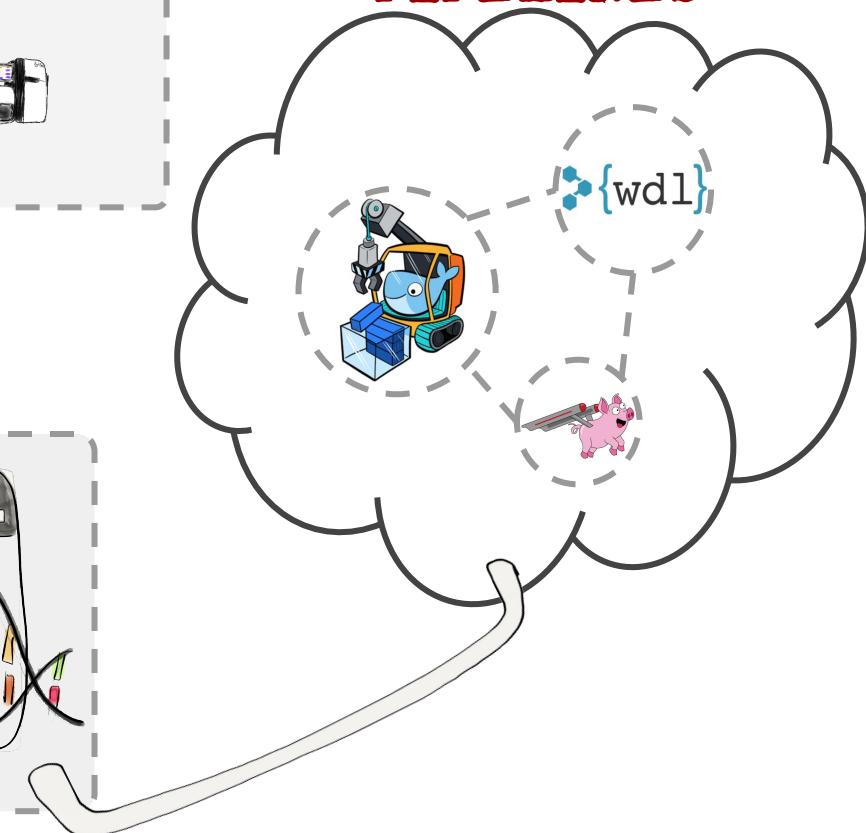
EXPERIMENT

Data model
Ontology
Audits (badges)
AB standards
File validation
Uniform Pipelines



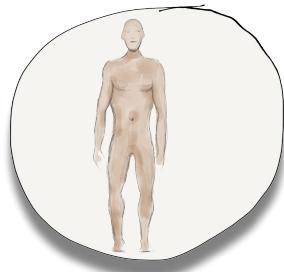
PORTAL

PIPELINES



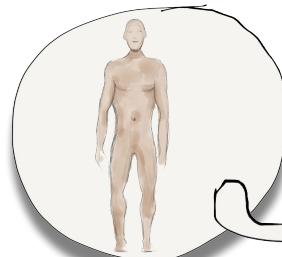
Data model

HumanDonor.json

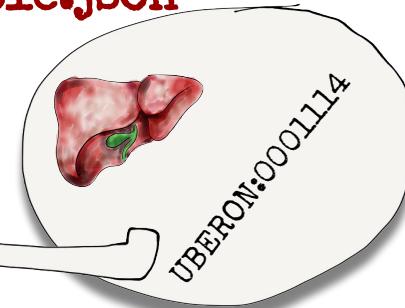


Data model
Ontology

HumanDonor.json



Biosample.json

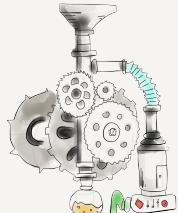


Experiment.json

Data model
Ontology
Audits (badges)

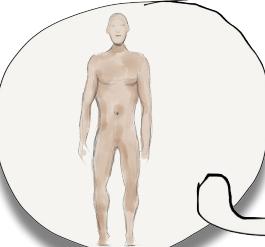


Library.json



CTCF ChIP-seq
OBI:0000716

HumanDonor.json



Biosample.json



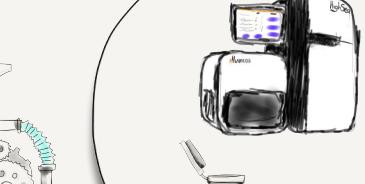
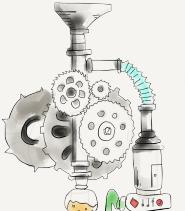
UBERON:0001114

Data model
Ontology
Audits (badges)
AB standards

Experiment.json

AntibodyLot.json

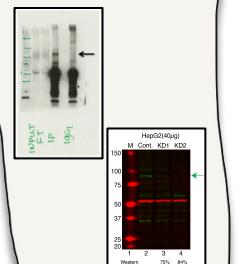
Library.json



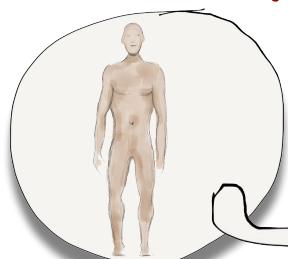
CTCF ChIP-seq
OBI:0000716



characterized to standards



HumanDonor.json



Biosample.json



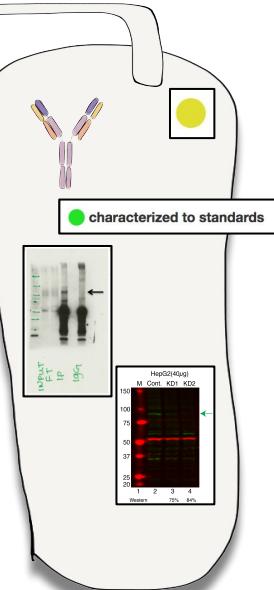
UBERON:0001114

Data model
Ontology
Audits (badges)
AB standards
File validation

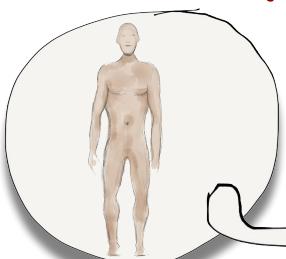
Experiment.json



AntibodyLot.json



HumanDonor.json

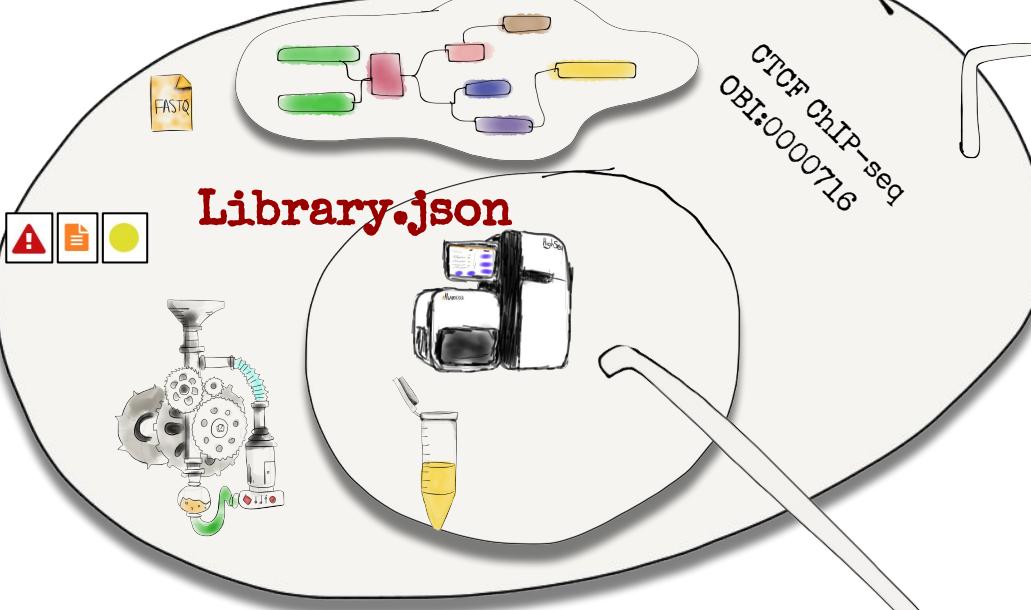


Biosample.json

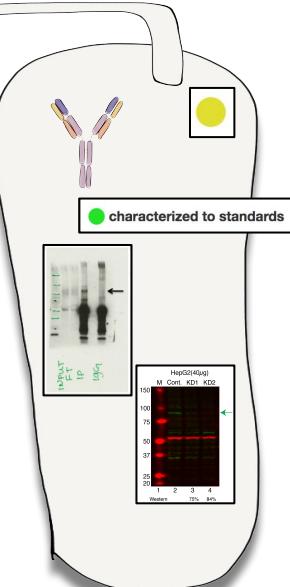


Data model
Ontology
Audits (badges)
AB standards
File validation
Uniform Pipelines

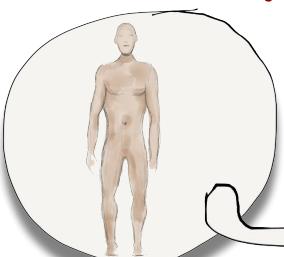
Experiment.json



AntibodyLot.json



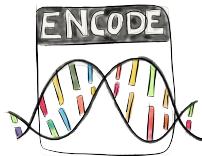
HumanDonor.json



Biosample.json



ENCODE Data Coordination Center



Data wranglers



Mike

Ben

Pipeline devs.



Jason

Esther

Idan

Seth

Otto

Bek

Tim

Zack

Kath

Nick

Keenan

Forrest

Casey

Karthik

Wrangler associates

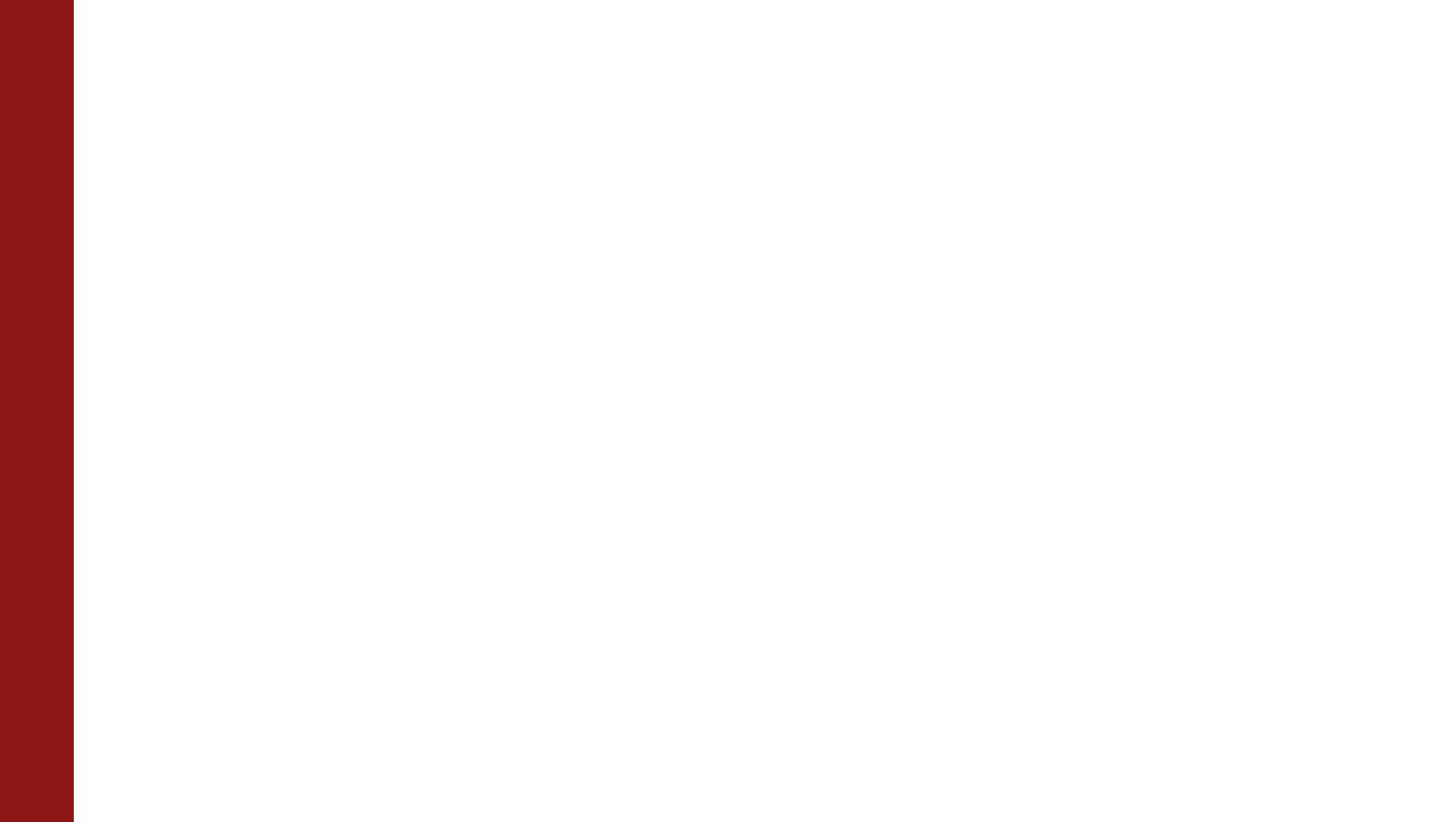
Software devs.

U24 HG009397 and U41 HG006992



National Human Genome
Research Institute





Challenges

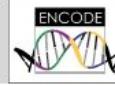
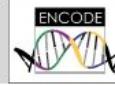
- protocol modeling
- standardization of emerging technologies:
 - genetic modification
 - functional characterization
- community data

Challenges

Genetic modification object

- method (RNAi, CRISPR, etc.)
- purpose (activation, repression, etc.)
- category (deletion, insertion, etc.)

Genetic modification

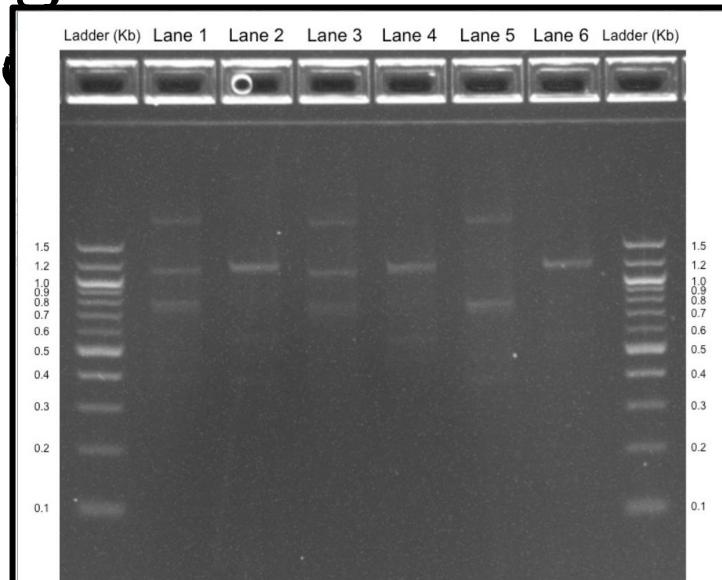
Summary		Attribution	
Description:	ChIP-Seq on HepG2	Lab:	Richard Myers, HAIB
Type:	insertion	Award PI:	Richard Myers, HAIB
Tags:	3xFLAG — C-terminal	Project:	ENCODE
Purpose:	tagging	Aliases:	richard-myers:SL272463-SL272464-gm
Modification site			
Target:	FLAG-NR5A1-human		
Modification method			
Technique:	CRISPR		
		 ENCODE PHASE 4	 ENCODE PHASE 4

Genetic modification

To the best of our knowledge there is no large scale systematic assessment of commercial modifications, ENCODE should? devote large effort towards conducting characteriations.

future challenges

Explain the challenge to standardize
the submitted characterizations



Lane 1: Homology arm 1 for Rep 1, expected size: 1105
Lane 2: Homology arm 2 for Rep 1, expected size: 933
Lane 3: Homology arm 1 for Rep 2, expected size: 1105
Lane 4: Homology arm 2 for Rep 2, expected size: 933
Lane 5: Homology arm 1 for wild type HepG2 background
Lane 6: Homology arm 2 for wild type HepG2 background