

Improving the value of public data with *recount2* and phenotype prediction



Shannon E. Ellis, PhD
Assistant Teaching Professor
Cognitive Science



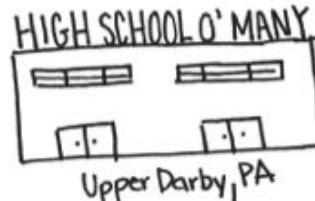
I LIKE
ALL
THE THINGS!



SHANNON



I LIKE
ALL
THE THINGS!

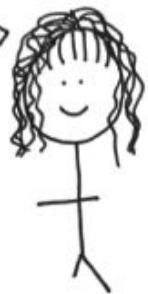


Upper Darby, PA

SCIENCE
IS PRETTY
COOL!

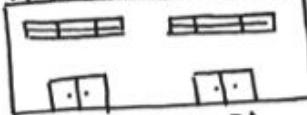


SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA



Wilkes-Barre, PA

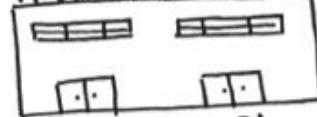
SCIENCE
IS PRETTY
COOL!

SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

GENETICS
IS AWESOME!



...
but what
did that
software
actually
do?

KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE
IS PRETTY
COOL!

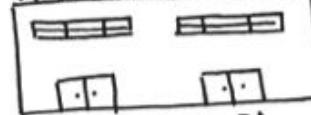


SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

JOHNS HOPKINS



Baltimore, MD

GENETICS
IS AWESOME!

...
but what
did that
software
actually
do?

KING'S COLLEGE



Wilkes-Barre, PA

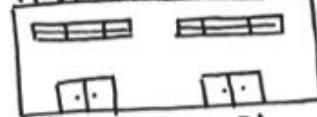
SCIENCE
IS PRETTY
COOL!

SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

JOHNS HOPKINS



Baltimore, MD

GENETICS
IS AWESOME!

but what
did that
software
actually
do?



KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE
IS PRETTY
COOL!



DATA
ANALYSIS
IS WHERE
IT'S AT!

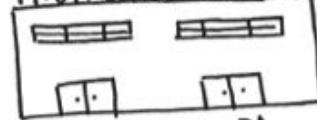


SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

JOHNS HOPKINS



Baltimore, MD

GENETICS
IS AWESOME!

but what
did that
software
actually
do?



KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE
IS PRETTY
COOL!



DATA
ANALYSIS
IS WHERE
IT'S AT!

JEFF
LEEK

COME DO
A POSTDOC
WITH ME!

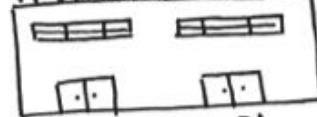


SHANNON



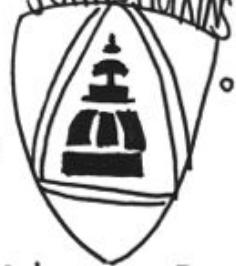
I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

JOHNS HOPKINS



Baltimore, MD

GENETICS
IS AWESOME!

but what
did that
software
actually
do?



KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE
IS PRETTY
COOL!



I LOVE
TEACHING!



JEFF
LEEK

COME DO
A POSTDOC
WITH ME!



DATA
ANALYSIS
IS WHERE
IT'S AT!

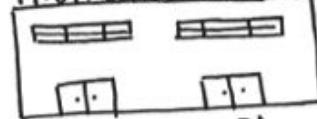


SHANNON



I LIKE
ALL
THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

JOHNS HOPKINS



Baltimore, MD

GENETICS
IS AWESOME!

but what
did that
software
actually
do?

KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE
IS PRETTY
COOL!



I LOVE
TEACHING!

TODAY

OPEN
SCIENCE



CHIANG
GENETICS

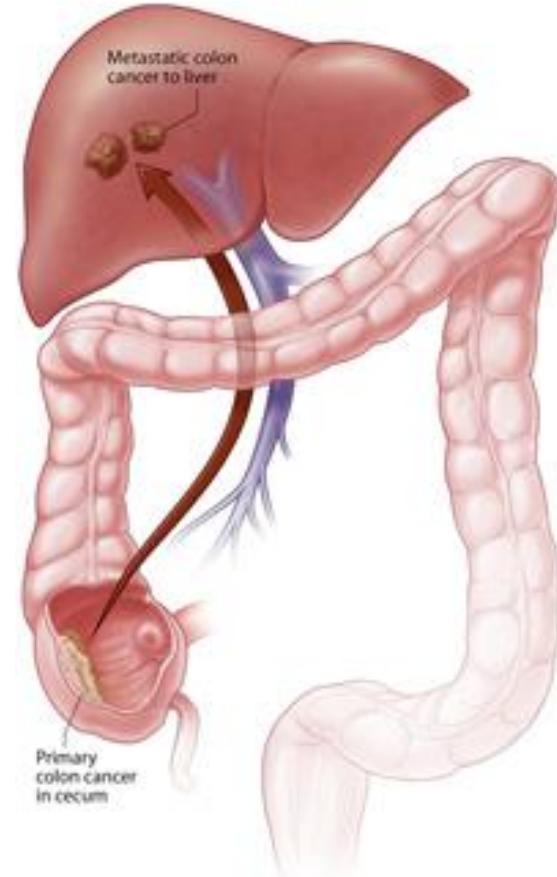


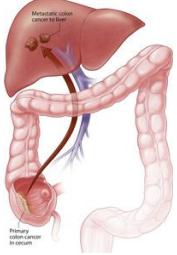
DATA
ANALYSIS
IS WHERE
IT'S AT!

JEFF
LEEK

COME DO
A POSTDOC
WITH ME!

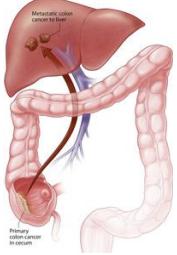
What makes primary cancer different than metastatic cancer?





What makes primary cancer different than metastatic cancer?

Find a
researcher
with access
to patient
samples

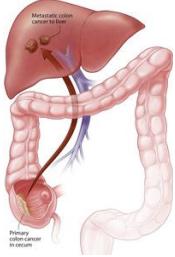


What makes primary cancer different than metastatic cancer?

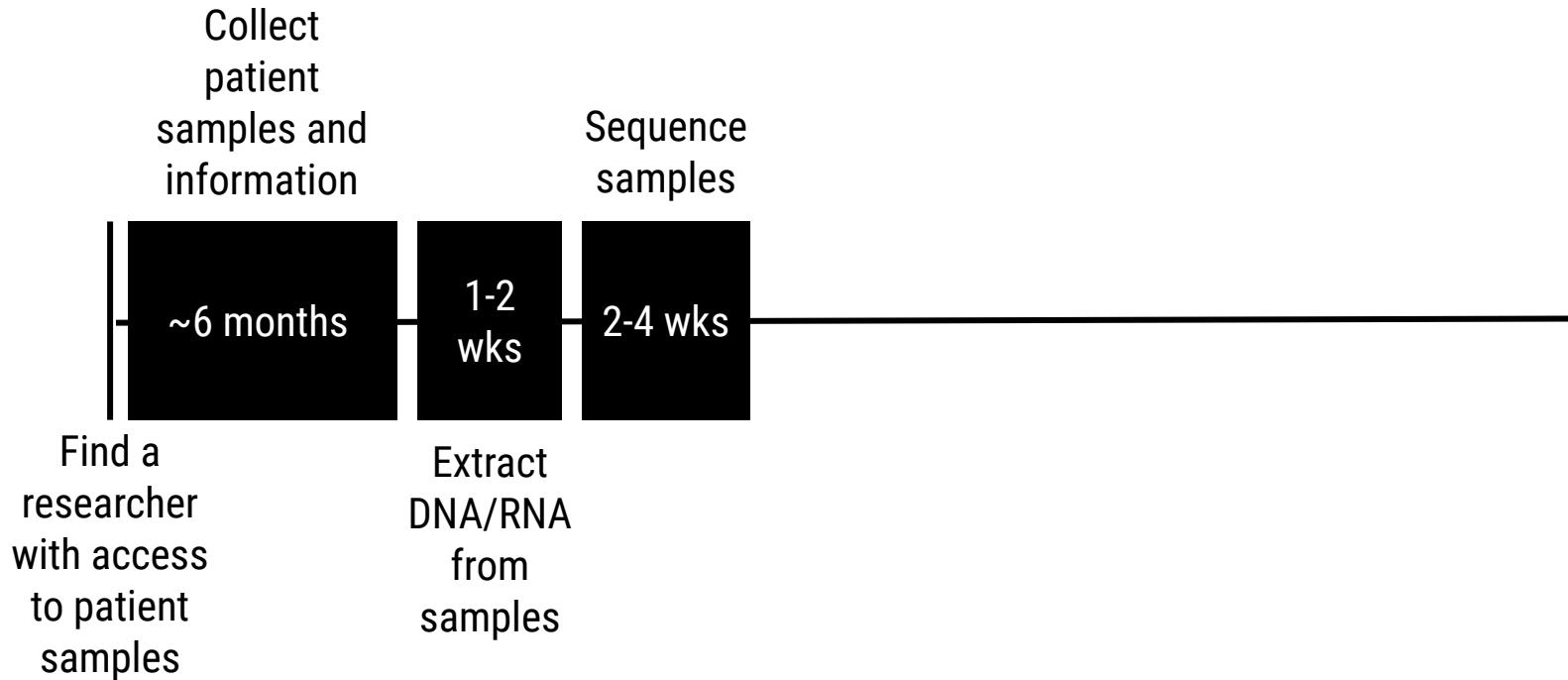
Collect
patient
samples and
information

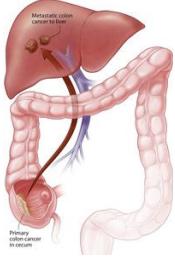
~6 months

Find a
researcher
with access
to patient
samples

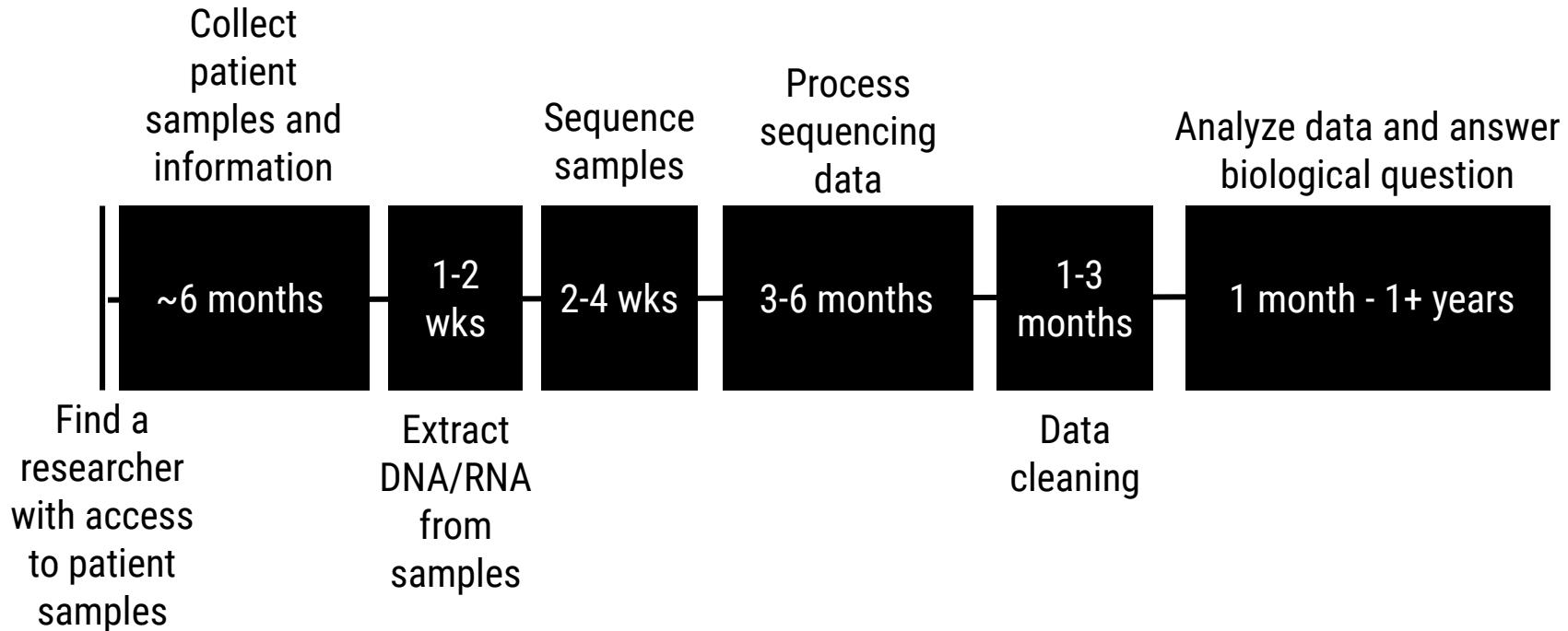


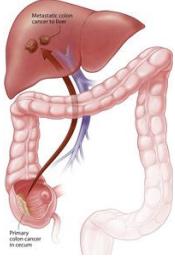
What makes primary cancer different than metastatic cancer?



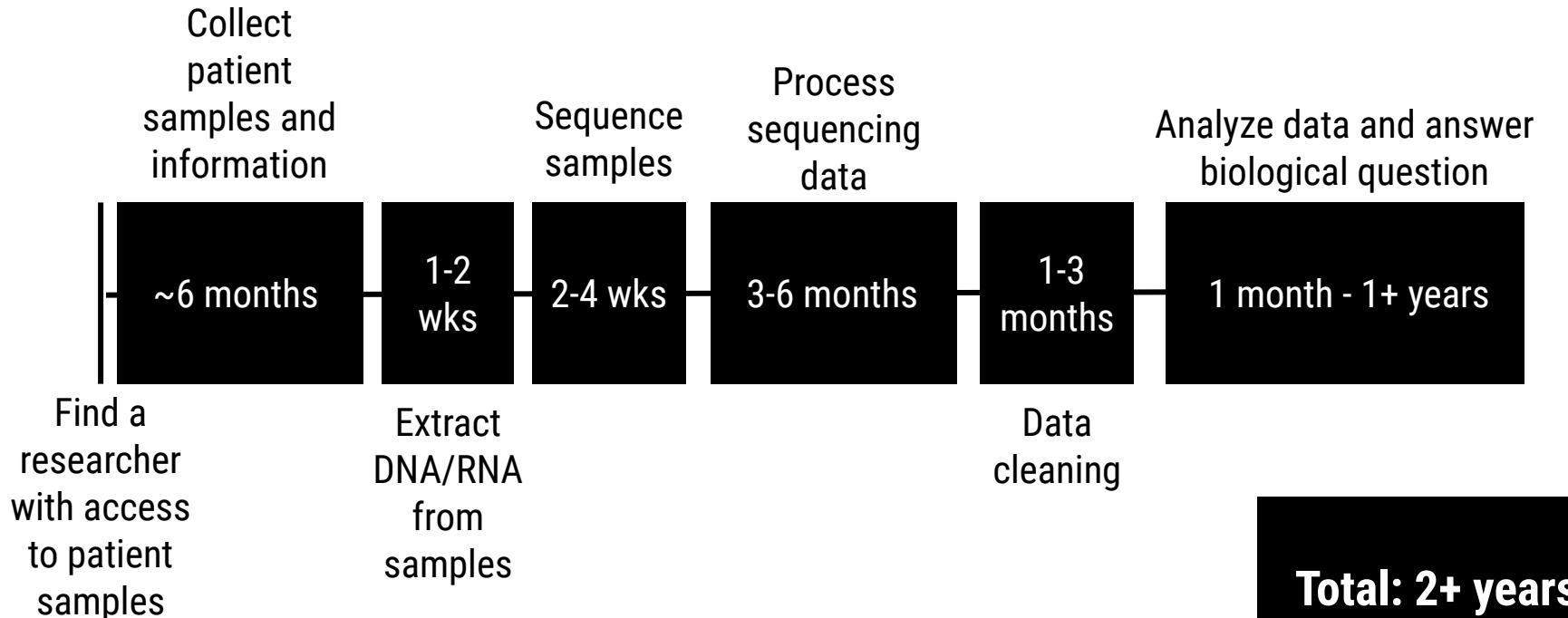


What makes primary cancer different than metastatic cancer?



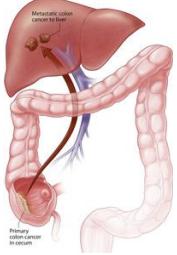


What makes primary cancer different than metastatic cancer?

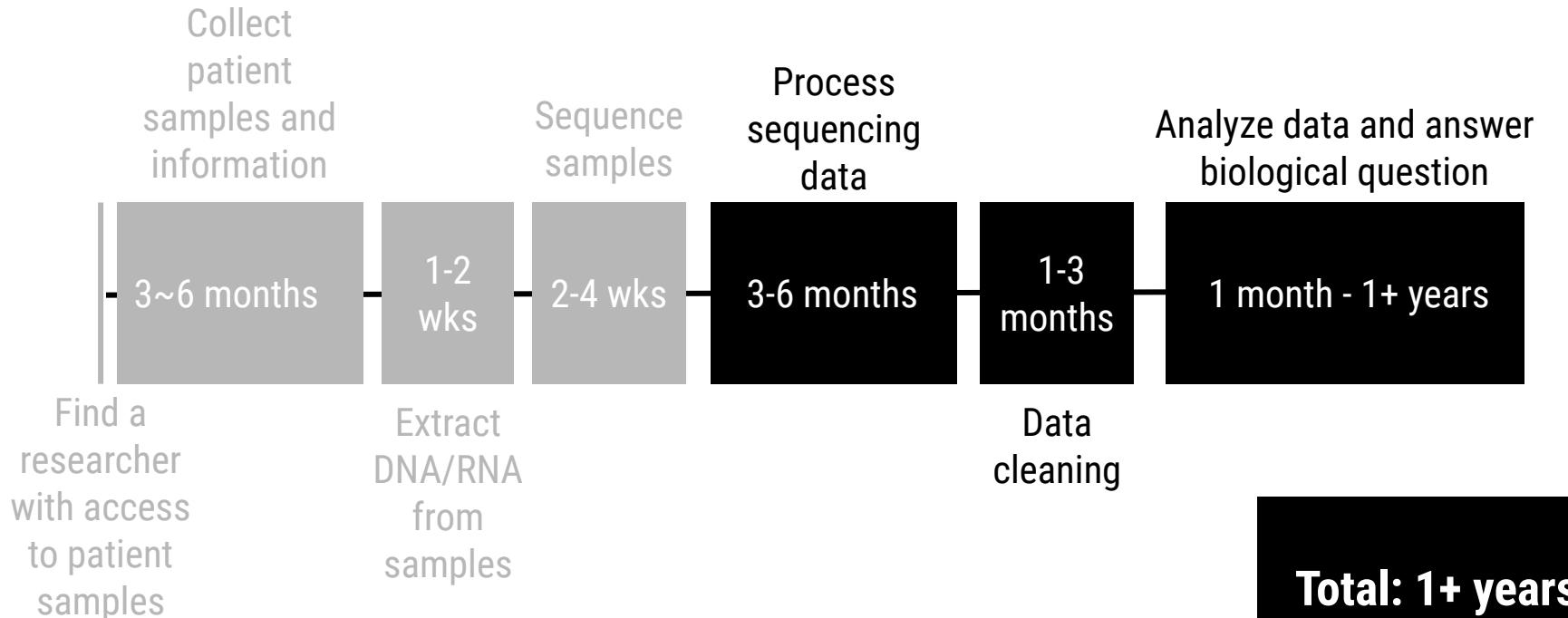




Biologists have recently gotten pretty good at making their data available to the public.



What makes primary cancer different than metastatic cancer?

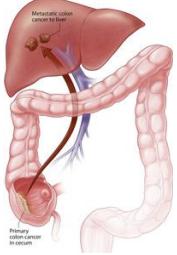




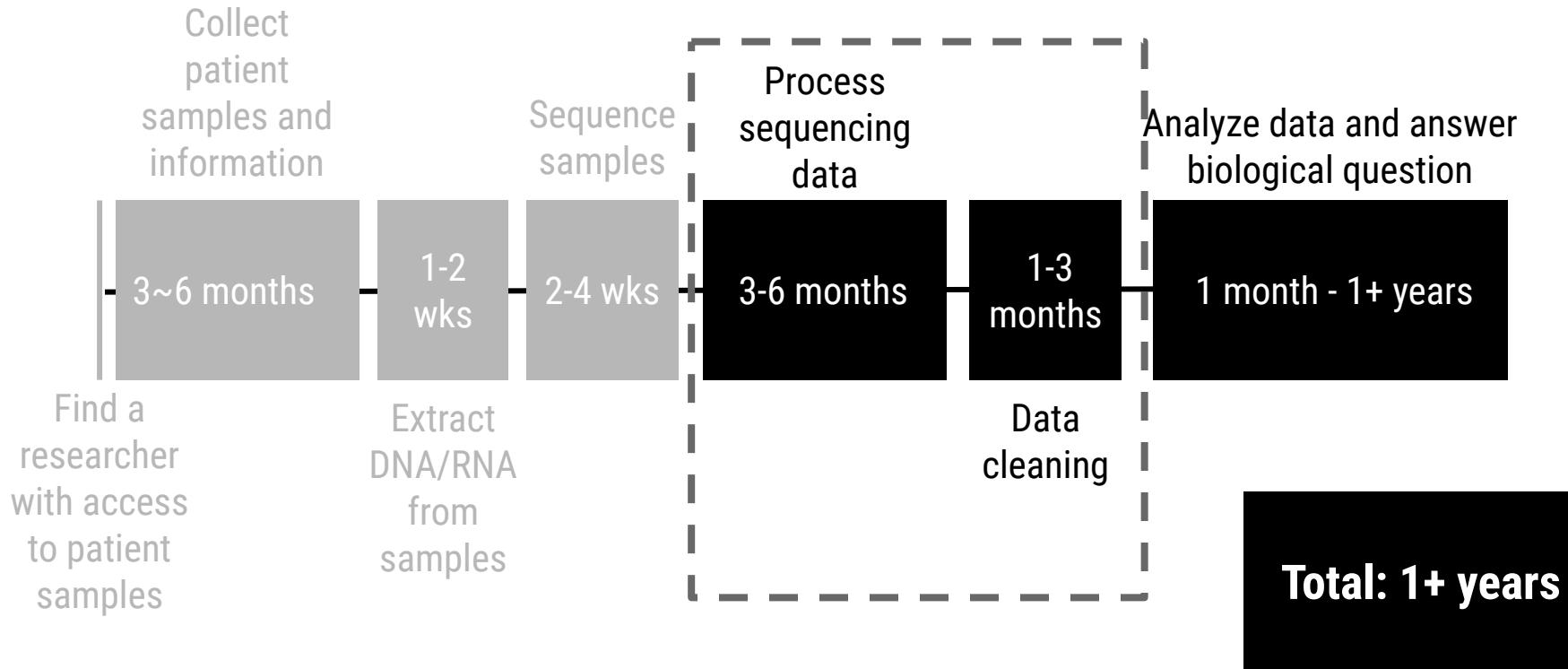
Biologists have recently gotten pretty good at making their data available to the public.



...but they're *not great* at making these data easily accessible and well-annotated.



What makes primary cancer different than metastatic cancer?



Project Goal: Take publicly available RNA-Seq data and make it available and easy-to-use

- Comprehensive
- Easy to Get
- Useful for future study

Genetics101



The Central Dogma of Genetics

DNA



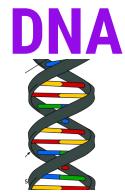
ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT
TGACTGGATCTAGTCAGCTAGCATATGCTAATGTTTAGTAGCCGTA

The Central Dogma of Genetics

DNA


ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT

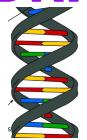
The Central Dogma of Genetics

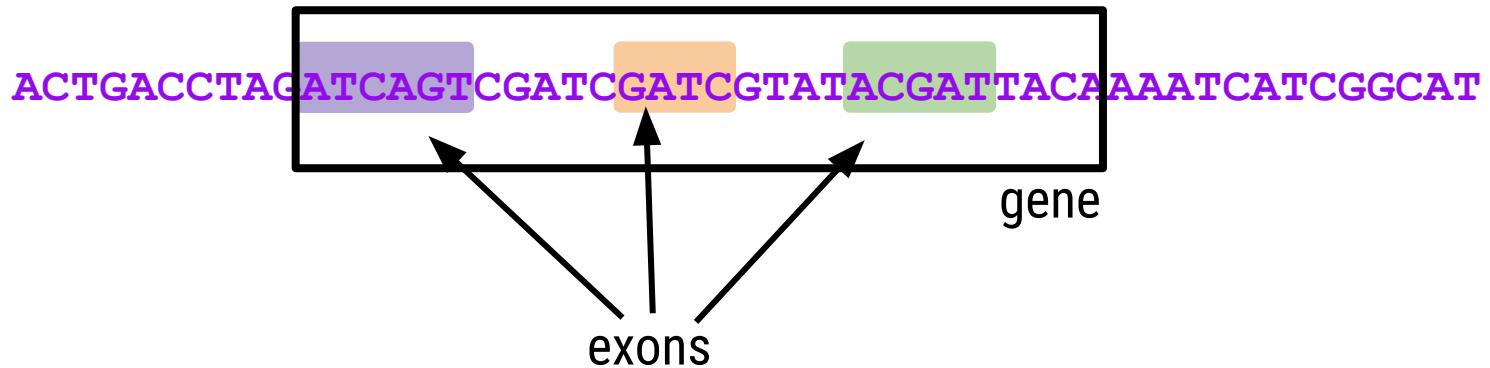


ACTGACCTAGATCAGTCGATCGATCGTATACGATTACA**AAATCATCGGCAT**

gene

The Central Dogma of Genetics

DNA




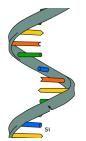
The Central Dogma of Genetics

DNA


ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT



transcription

RNA


AUCAGUCGAUCACCGAU

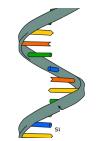
The Central Dogma of Genetics

DNA


ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT



transcription

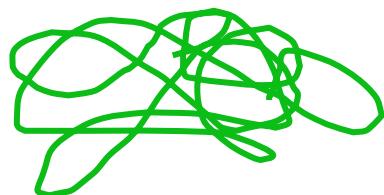
RNA


AUCAGUCGAUCACCGAU

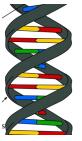
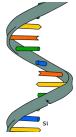


translation

proteins

Two copies of DNA -> many transcripts -> many proteins

<i>role in the cell</i>	<i># copies/cell</i>	<i>functional unit</i>	<i># unique functional units</i>
DNA 	<i>blueprint</i>	2	<i>gene</i>
RNA 			
proteins 			

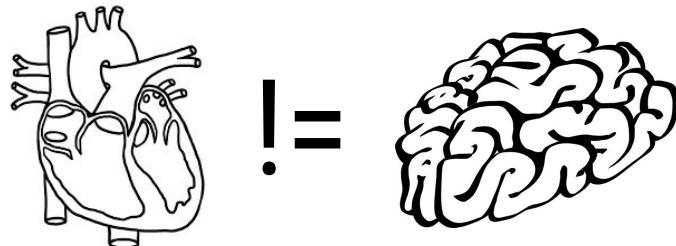
Two copies of DNA -> many transcripts -> many proteins

	<i>role in the cell</i>	<i># copies/cell</i>	<i>functional unit</i>	<i># unique functional units</i>
DNA	<i>blueprint</i>	2	<i>gene</i>	20,000
RNA				
proteins	<i>carry out cellular functions</i>	<i>varies $\sim 10^{10}$</i>	<i>proteins (metabolites, hormones, etc.)</i>	<i>~100,000</i>

Two copies of DNA -> many transcripts -> many proteins

	<i>role in the cell</i>	<i># copies/cell</i>	<i>functional unit</i>	<i># unique functional units</i>
DNA	<i>blueprint</i>	2	<i>gene</i>	20,000
RNA	<i>messenger</i>	<i>varies ~360,000</i>	<i>transcript</i>	<i>~100,000</i>
proteins	<i>carry out cellular functions</i>	<i>varies ~10¹⁰</i>	<i>proteins (metabolites, hormones, etc.)</i>	<i>~100,000</i>

Variability at the level of RNA
allows for a heart cell to
function differently than a
brain cell

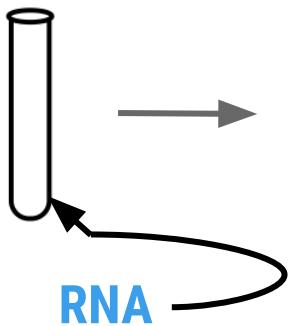


Measuring RNA levels



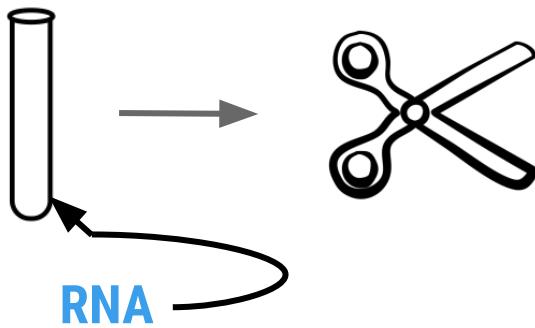
Next Generation Sequencing (NGS)
Has Completely Revolutionized
How We Study Genetics

Next Generation Sequencing (NGS) in one slide



Step 1: Extract
RNA to get
sample of
interest

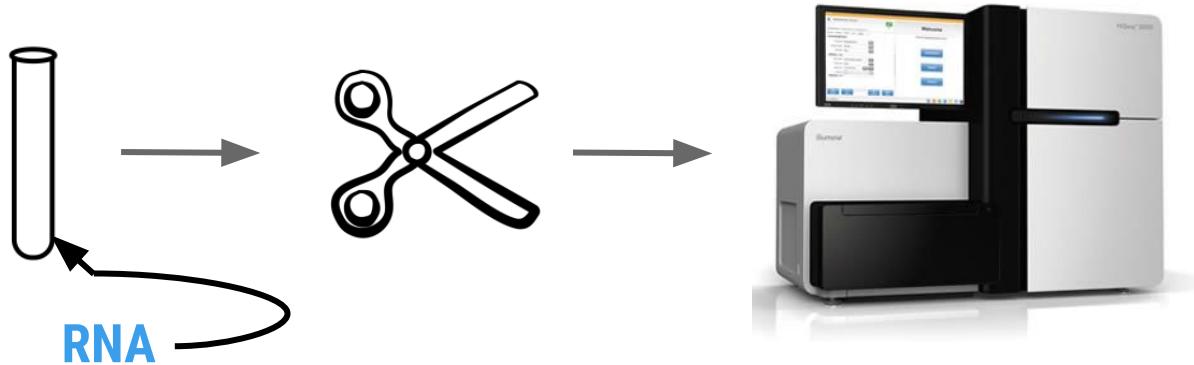
Next Generation Sequencing (NGS) in one slide



Step 1: Extract RNA to get sample of interest

Step 2: Chop up RNA into smaller pieces

Next Generation Sequencing (NGS) in one slide



Step 1: Extract RNA to get sample of interest

Step 2: Chop up RNA into smaller pieces

Step 3: Sequence the sample

Next Generation Sequencing (NGS) in one slide



Step 1: Extract RNA to get sample of interest

Step 2: Chop up RNA into smaller pieces

Step 3: Sequence the sample

Step 3: Obtain short read data from the sequencer

Next Generation Sequencing (NGS) in one slide



A short read tells you the sequence
of the RNA in that read

@22:16362385-16362561W:ENST00000440999:2:177:-40:244:S/2
CCAGCCCACCTGAGGCTTCTTTCTTCCAAGCCACATCACCACCTGGTGGAACTCTC
CTGTGAGGACAGCCA
+
GGFF<BB=>GBGIIIIIIIIIIIIIEGEHGHIIIIIIHFHB2/:=?EGGGEGFHH
IHEDBD?@@DDHHD
@22:16362385-16362561W:ENST00000440999:3:177:-56:294:S/2
GCGTGAGGCCACAGGGCCCAGCCCACCTGAGGCTTCTTTCTTCCAAGCCACATCACCA
TCCTGGTGGAACTCT
+
@=ABBBBIIIIIIIIHHGGGGIIDBDIIIIIIIGIIIIHIIIIHFDD@BBDBGGFIDEE8DC
C/29>BGFCGHHGF
@22:16362385-16362561W:ENST00000440999:4:177:137:254:S/1
TCACCATCCTGGTGGAACTCTCCTGTGAGGACAGCCAAGGCCTGAACCTACCTGCaGTGGGG
AGCACCTCAGGGTTT
+
DDGBBCGGGIGGGBDDDHIIGGDGD77=BDIIIIIIIFHHHIIIIHEFHFDD@>DEC
HHIFDDHH8@BEDDI
@22:16362385-16362561W:ENST00000440999:5:17:18:15:S/2
AGGGTTGCCAGGCAACCAGCCAGCCCTGGTCCAAGGCATCCTGGAGCGAGTTGTGGATG
GCAAAAAGACNCGCC
+
HIGHIHFEHE4111:.;8@?@HDIIIIIIIEGGIHHIIIGA?=:FIIIDD8.02506A8=

40M+



Sequence Identifier

@22:16362385-16362561W:ENST00000440999:2:177:-40:244:S/2

Sequence

CCAGCCCACCTGAGGCTTCTTTCTTCCAAGCCACATCACCATCCTGGTGGAACTCTC
CTGTGAGGACAGCCA

+

GGFF<BB=>GBGIIIIIIIIIIIEGEHGHIIIIIIHFBB2/:=?EGGEGFHH
IHEDBD?@@DDHHD

@22:16362385-16362561W:ENST00000440999:3:177:-56:294:S/2

GCGTGAGGCCACAGGGCCCAGCCCACCTGAGGCTTCTTTCTTCCAAGCCACATCACCA
TCCTGGTGGAACTCT

+

@=ABBBIIIIIIIIHHGGGGIIDBDIIIIIGIIIIHIIIIHFDD@BBDBGGFIDEE8DC
C/29>BGFCGHHGF

@22:16362385-16362561W:ENST00000440999:4:177:137:254:S/1

TCACCATCCTGGTGGAACTCTCCTGTGAGGACAGCCAAGGCCTGAACCTACCTGCaGTGGGG
AGCACCTCAGGGTTT

+

DDGBBCGGGIGGGBDDDHIIGGDGD77=BDIIIIIIIFHHHIIIIHEFHFDD@>DEC
HHIFDDHH8@BEDDI

@22:16362385-16362561W:ENST00000440999:5:1,7:18,15:S/2

AGGGTTGCCAGGCAACCAGCCAGCCCTGGTCCAAGGCATCCTGGAGCGAGTTGTGGATG
GCAAAAAGACNCGCC

+

HIGHIHFEHE4111:.;8@?@HDIIIIIIIEGGIHHIIIGA?=:FIIIDD8.02506A8=

40M+



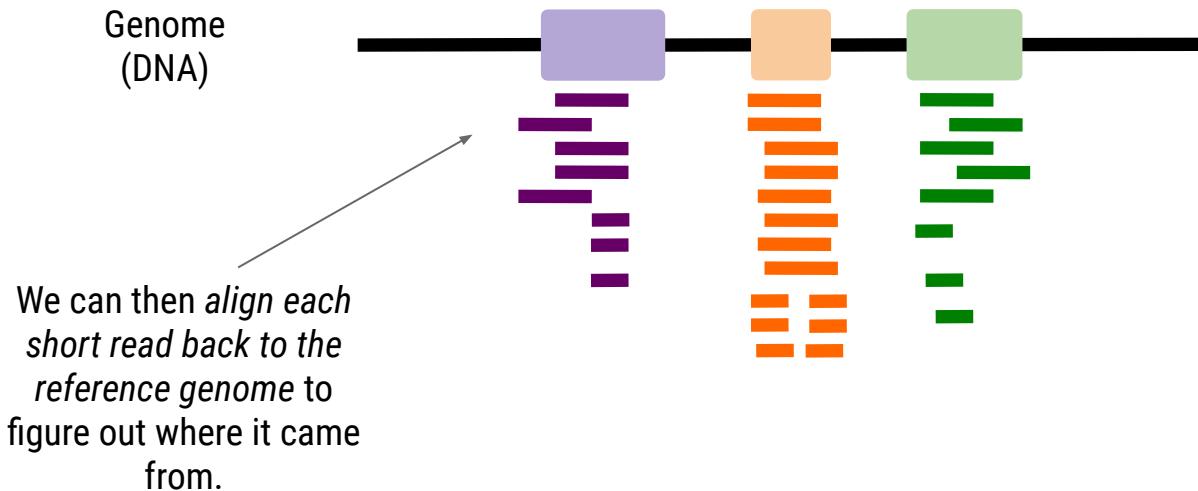
We've got 40M+ reads.
What does that all mean?

We first need to align these reads back to the genome



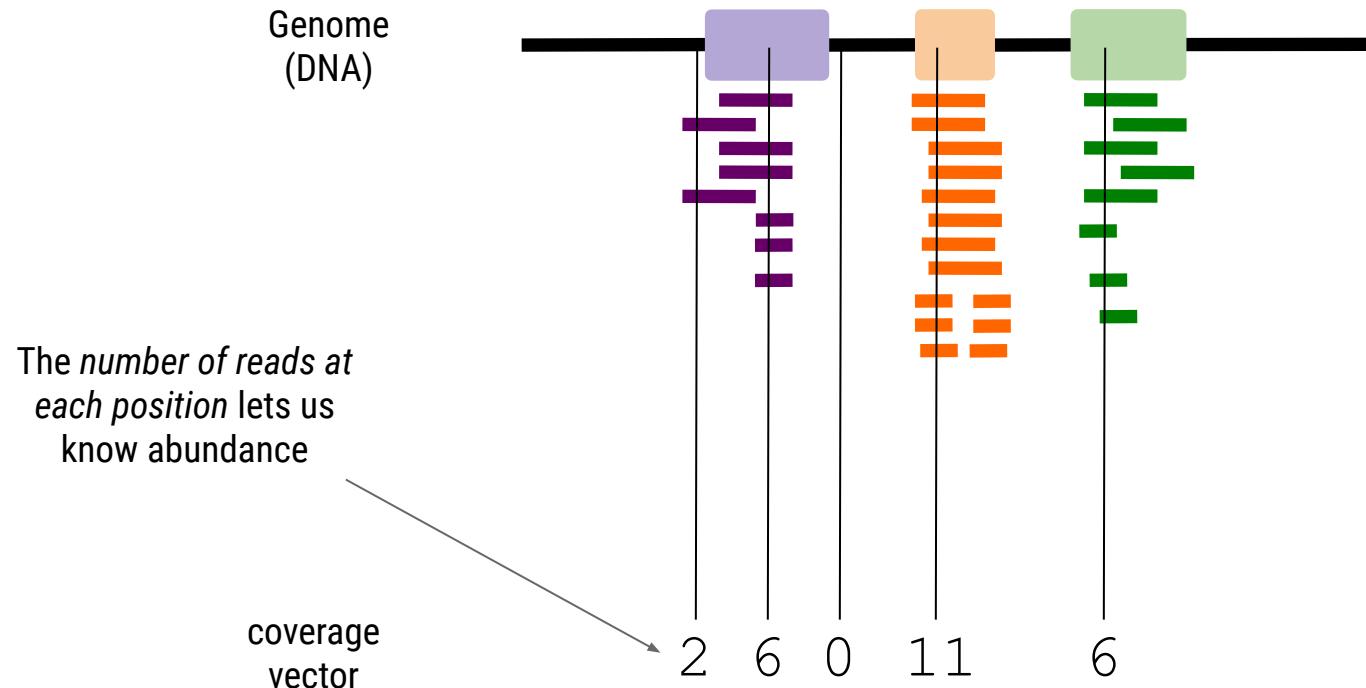
The Human Genome Project (HGP) determined the *reference sequence* of the human genome in 2001.

We first need to align these reads back to the genome

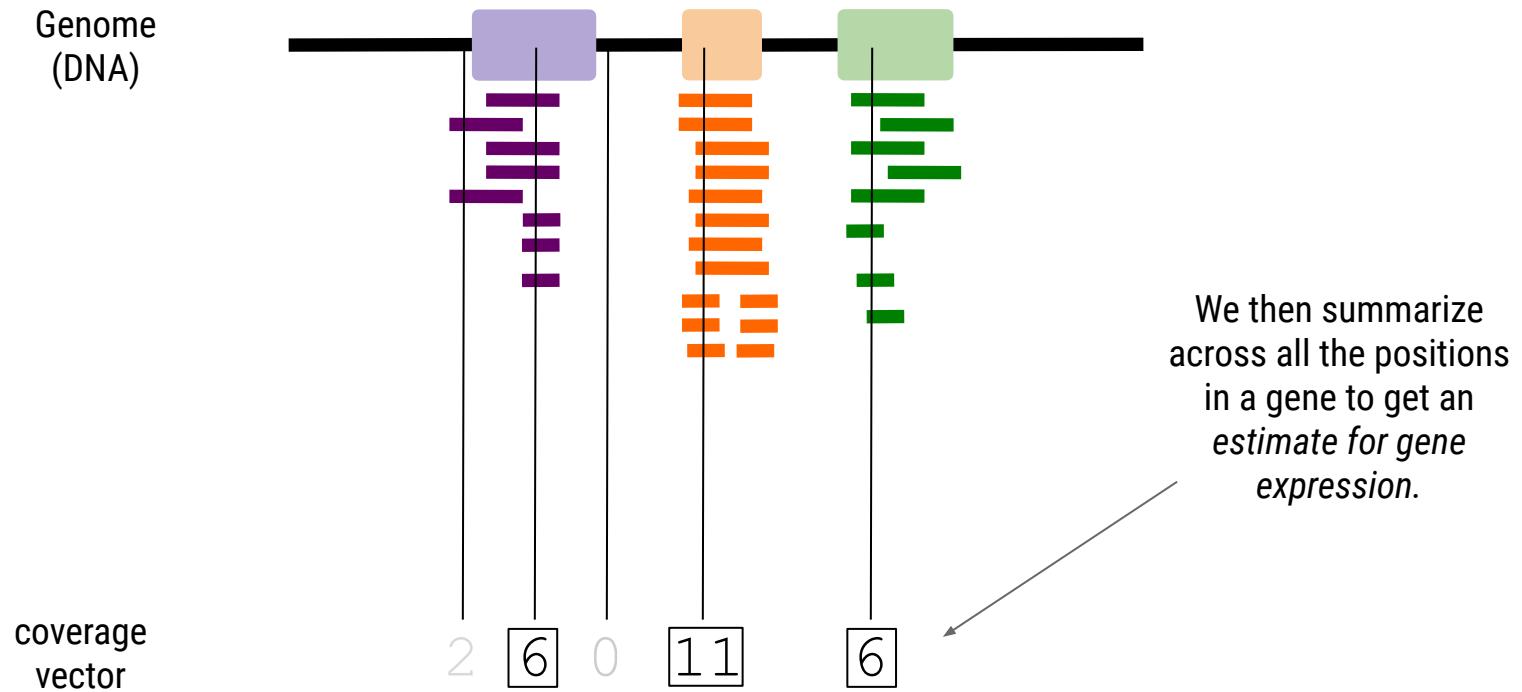


Nellore et al. (2016)
Bioinformatics
<http://rail.bio/>

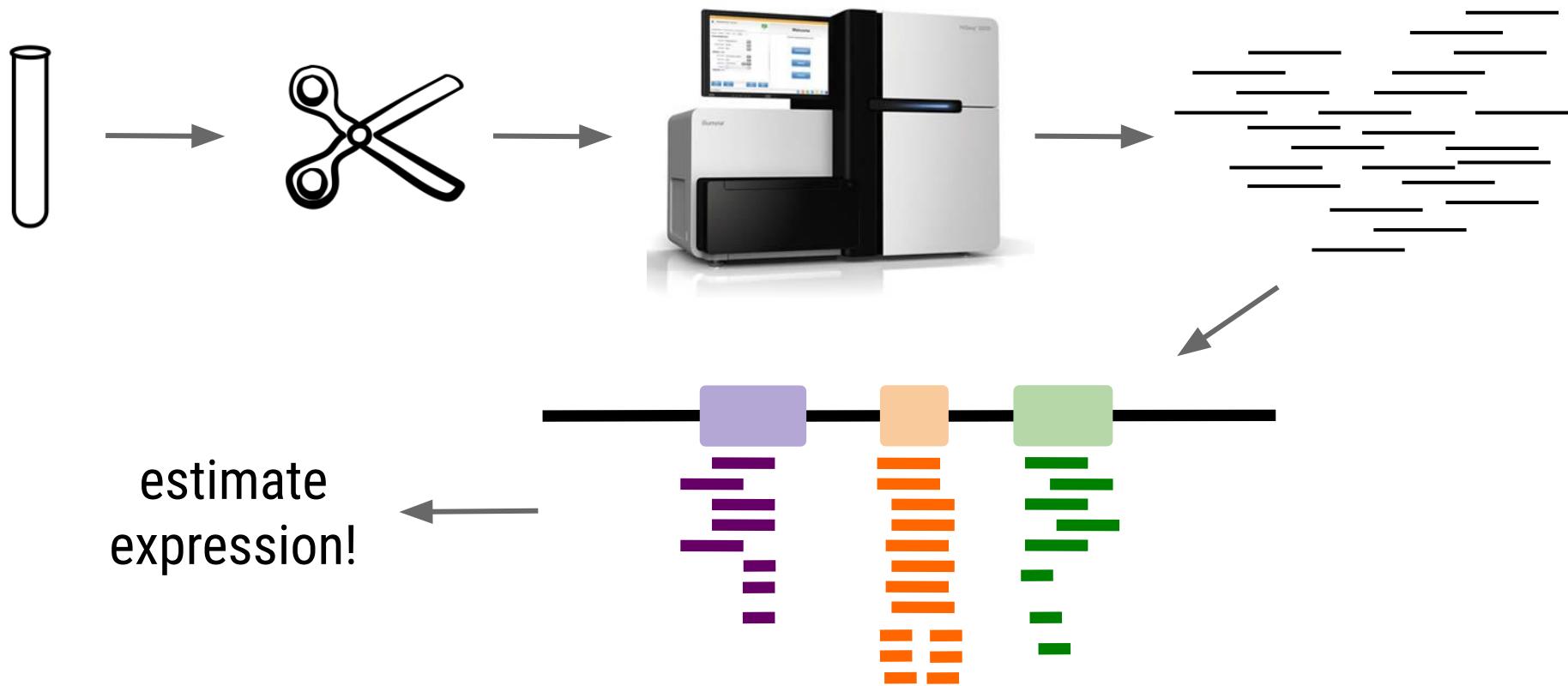
We first need to align these reads back to the genome



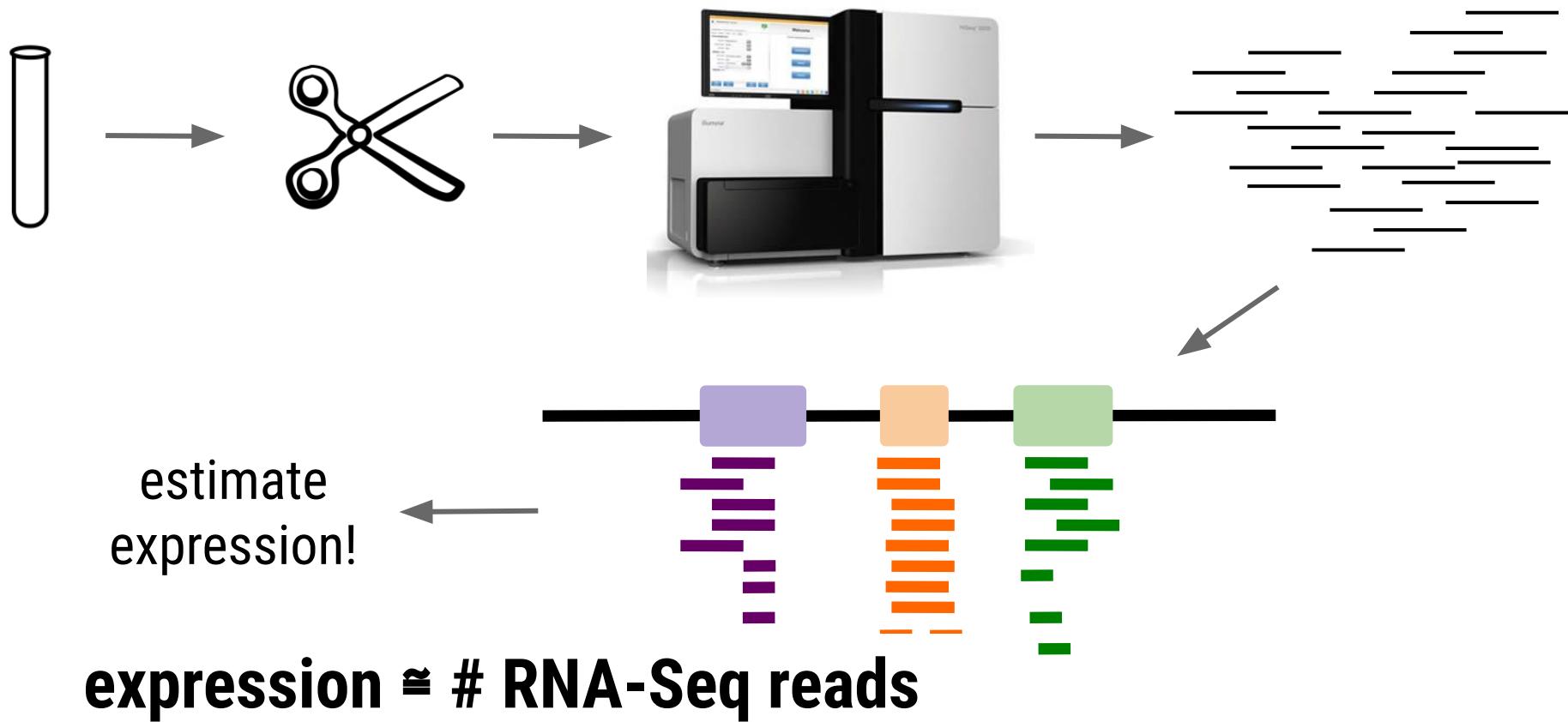
We first need to align these reads back to the genome



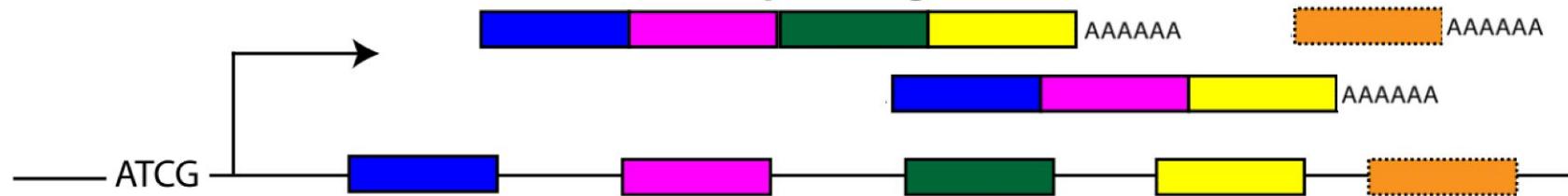
RNA-Seq = estimate expression across entire genome



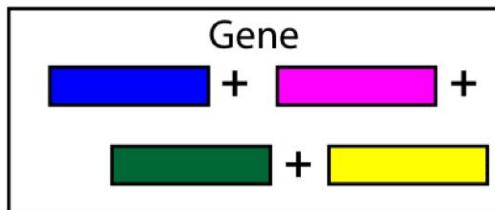
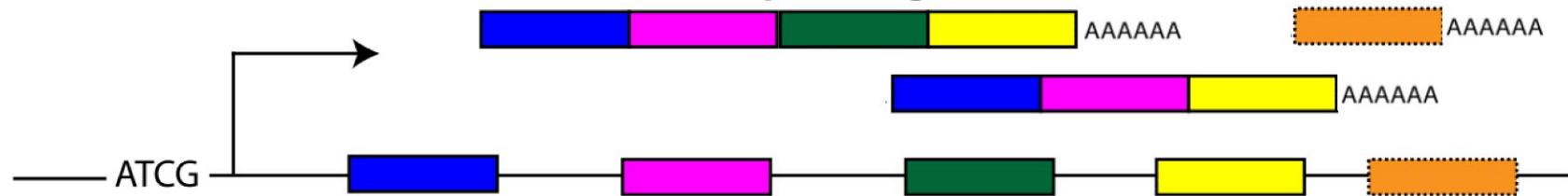
RNA-Seq = estimate expression across entire genome



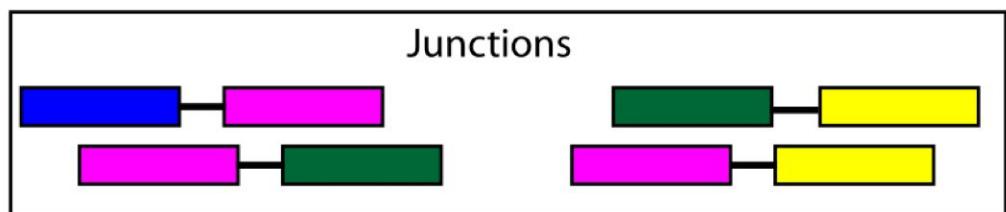
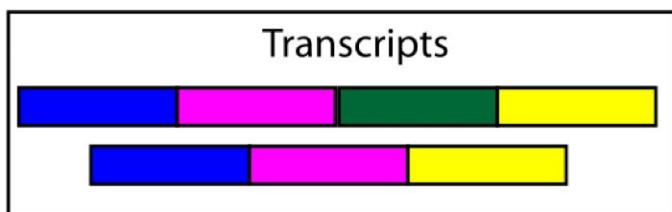
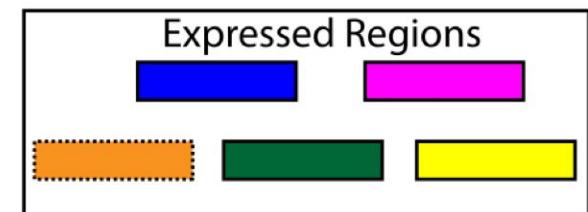
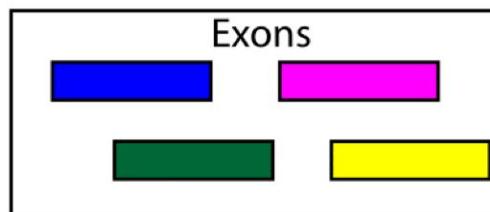
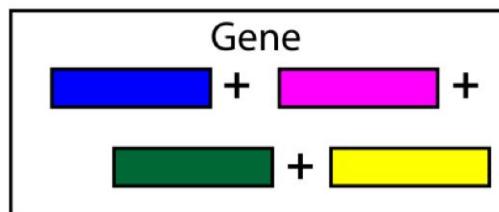
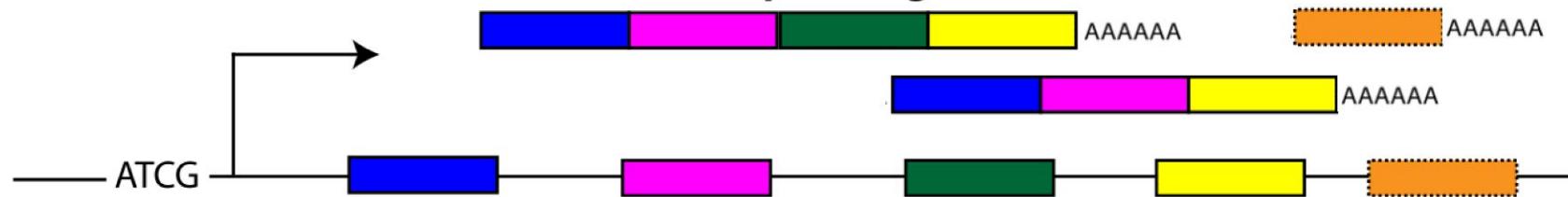
RNA Sequencing



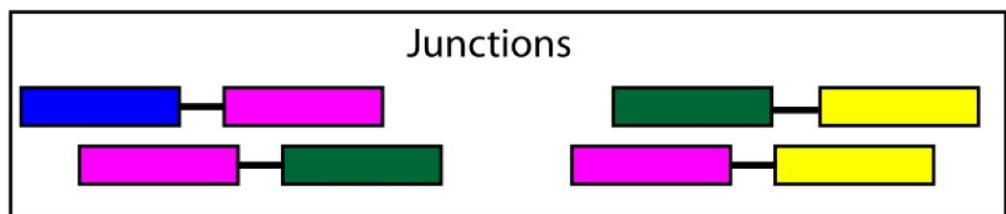
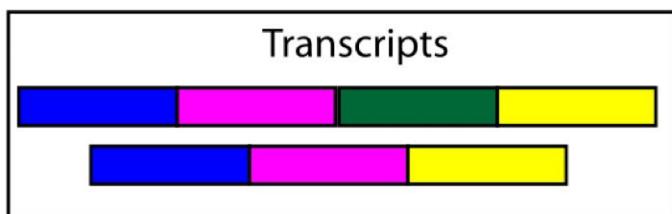
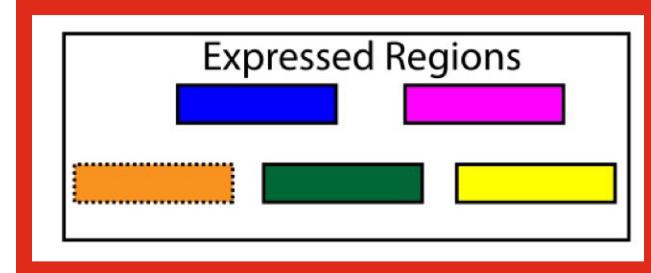
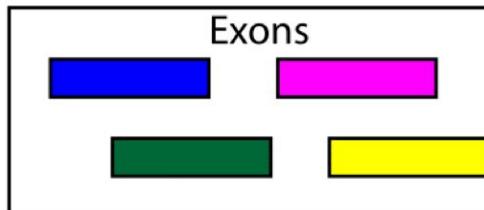
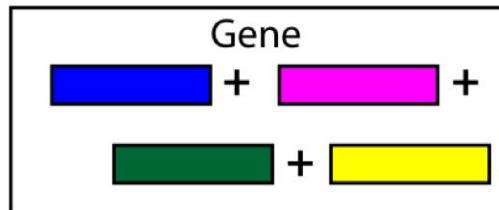
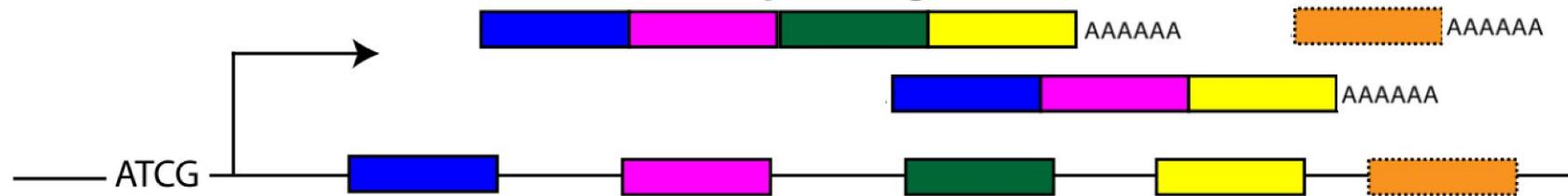
RNA Sequencing

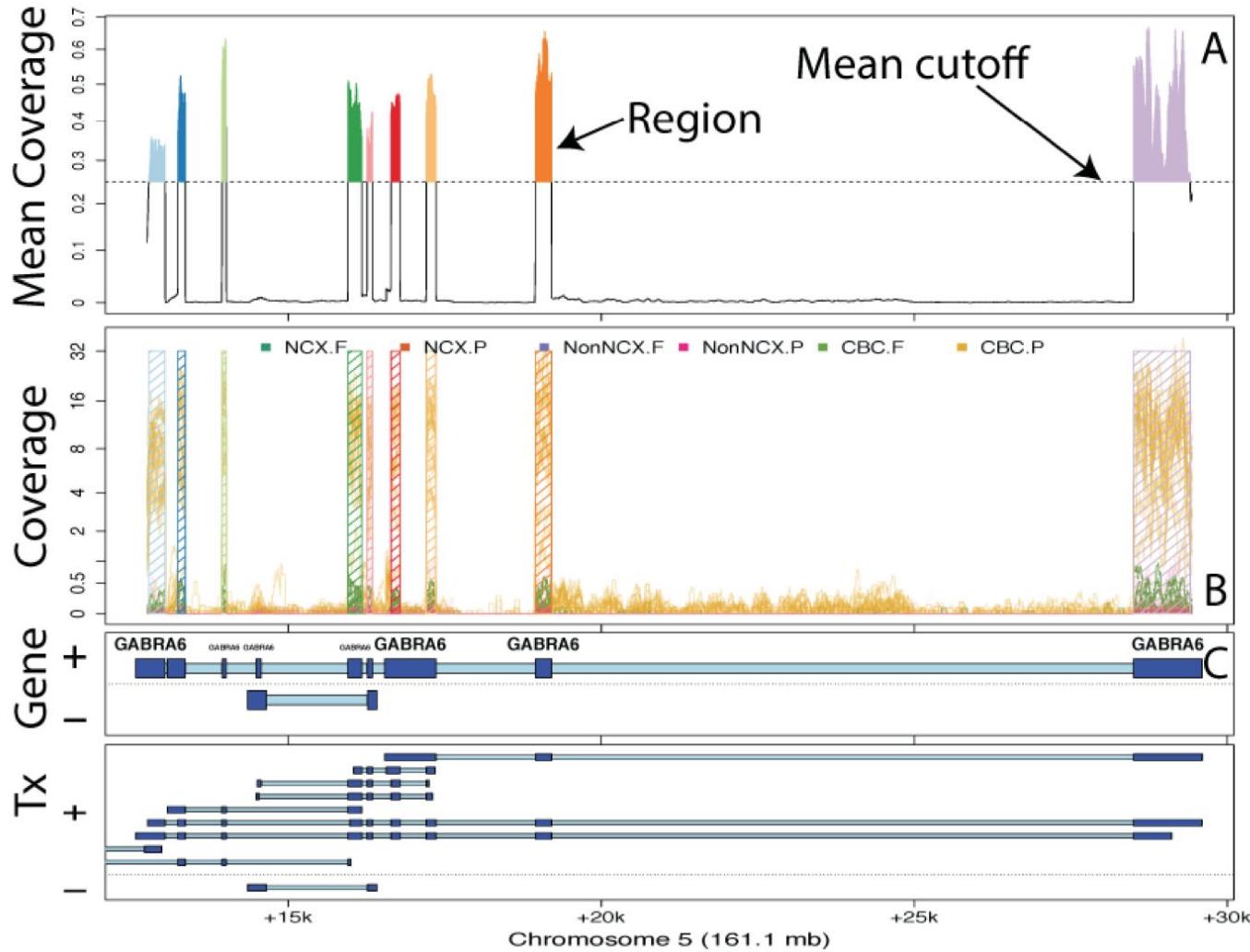


RNA Sequencing



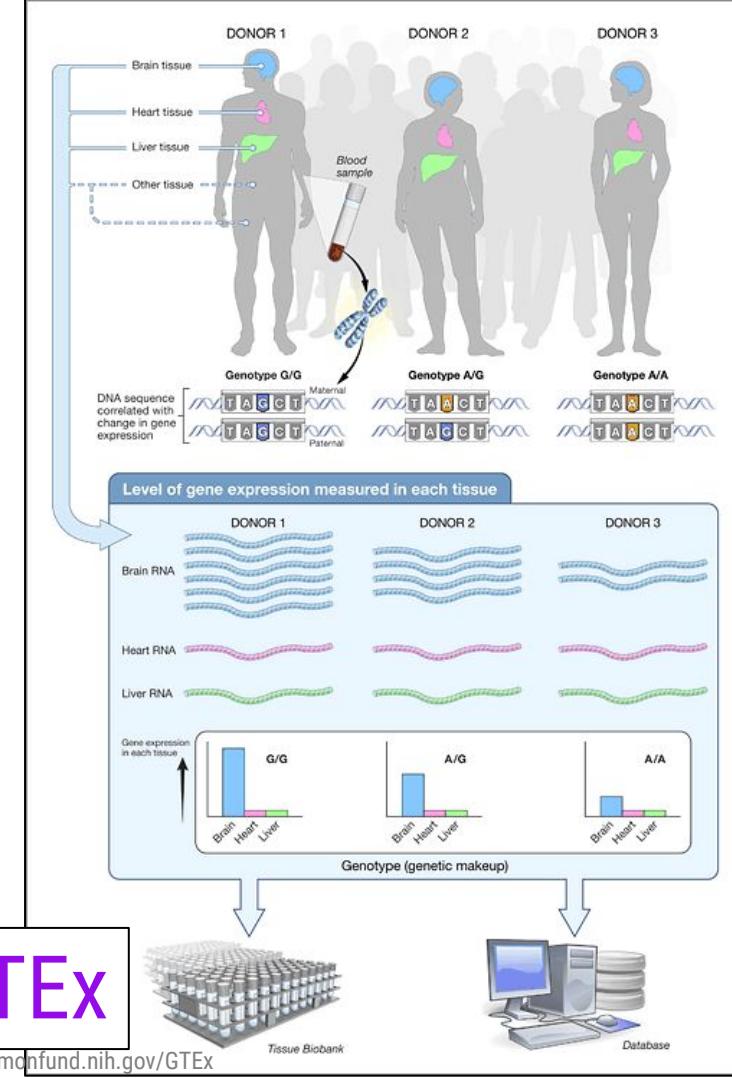
RNA Sequencing

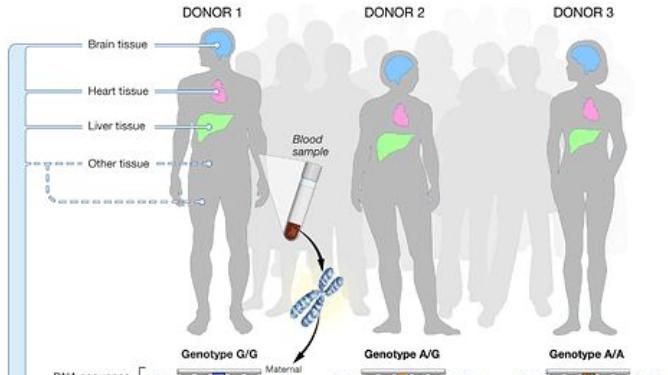




Project Goal: Take publicly available RNA-Seq data and make it available and easy-to-use

- Comprehensive
- Easy to Get
- Useful for future study





GTEx

<https://commonfund.nih.gov/GTEx>



NATIONAL CANCER INSTITUTE THE CANCER GENOME ATLAS

TCGA BY THE NUMBERS

TCGA produced over

2.5 PETABYTES
of data

To put this into perspective, 1 petabyte of data is equal to

212,000 DVDs



TCGA data describes

33 DIFFERENT TUMOR TYPES
10 RARE CANCERS

...based on paired tumor and normal tissue sets collected from

11,000 PATIENTS

...using
7 DIFFERENT DATA TYPES



TCGA RESULTS & FINDINGS

MOLECULAR BASIS OF CANCER
Improved our understanding of the genomic underpinnings of cancer

TUMOR SUBTYPES
Revolutionized how cancer is classified

THERAPEUTIC TARGETS
Identified genomic characteristics of tumors that can be targeted with currently available therapies or used to help with drug development

For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the serous subtype of ovarian cancer on a molecular level, suggesting that despite arising from different tissues in the body, these subtypes may share a common path of development and respond to similar therapeutic strategies.

TCGA revolutionized how cancer is classified by identifying tumor subtypes with distinct sets of genomic alterations.*

TCGA's identification of targetable genomic alterations in lung squamous cell carcinoma led to NCI's Lung-MAP Trial, which will treat patients based on the specific genomic changes in their tumor.

THE TEAM

20 COLLABORATING INSTITUTIONS
across the United States and Canada



WHAT'S NEXT?

The Genomic Data Commons (GDC) houses TCGA and other NCI-generated data sets for scientists to access from anywhere. The GDC also has many expanded capabilities that will allow researchers to answer more clinically relevant questions with increased ease.

TCGA

Analysis of stomach cancer revealed that it is not a single disease, but a disease composed of subtypes, including a new subtype characterized by infection with Epstein-Barr virus.

www.cancer.gov/ccg

SRA

SRA

 Advanced [Search](#) [Help](#)

SRA

Sequence Read Archive (SRA) makes biological sequence data available to the research community to enhance reproducibility and allow for new discoveries by comparing data sets. The SRA stores raw sequencing data and alignment information from high-throughput sequencing platforms, including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD System®, Helicos Heliscope®, Complete Genomics®, and Pacific Biosciences SMRT®.

Getting Started

[Understanding and Using SRA](#)[How to Submit](#)[Login to Submit](#)[Download Guide](#)

Tools and Software

[Download SRA Toolkit](#)[SRA Toolkit Documentation](#)[SRA-BLAST](#)[SRA Run Browser](#)[SRA Run Selector](#)

Related Resources

[dbGaP Home](#)[Trace Archive Home](#)[BioSample](#)[GenBank Home](#)

SRA

Project	No. of Sample
GTEx Genotype-Tissue Expression Project	9,962
TCGA The Cancer Genome Atlas	11,284
SRA Sequence Read Archive	49,848

Project	No. of Sample
GTEx Genotype-Tissue Expression Project	9,962
TCGA The Cancer Genome Atlas	11,284
SRA Sequence Read Archive	49,848

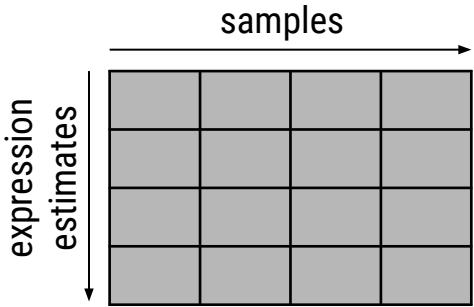
We'll take these ~70,000 samples, align each back to
the reference genome, and then, for each sample,
we'll estimate expression across the genome.

Project Goal: Take publicly available RNA-Seq data and make it available and easy-to-use

- Comprehensive
- Easy to Get
- Useful for future study

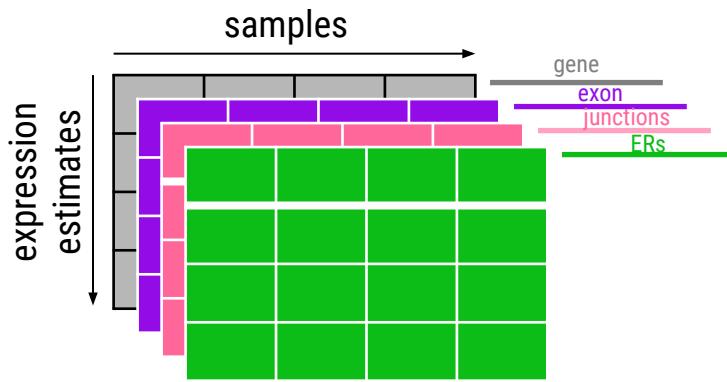


expression data for ~70,000 human samples





expression data for ~70,000 human samples



recount2: analysis-ready RNA-seq gene and exon counts datasets

[Datasets](#)[Popular datasets](#)[GTEx](#)[TCGA](#)[Documentation](#)[Download data with R](#)[Accessing recount2 via SciServer](#)[Contribute your data](#)

A multi-experiment resource of analysis-ready RNA-seq gene and exon count datasets

reCount2 is an online resource consisting of RNA-seq gene and exon counts as well as coverage bigWig files for 2041 different studies. It is the second generation of the [ReCount project](#). The raw sequencing data were processed with [Rail-RNA](#) as described at [bioRxiv 038224](#) which created the coverage bigWig files. For ease of statistical analysis, for each study we created count tables at the gene and exon levels and extracted phenotype data, which we provide in their raw formats as well as in RangedSummarizedExperiment R objects (described in the [SummarizedExperiment](#) Bioconductor package). We also computed the mean coverage per study and provide it in a bigWig file, which can be used with the [derfinder](#) Bioconductor package to perform annotation-agnostic differential expression analysis at the expressed regions-level as described at [bioRxiv 015370](#). The count tables, RangedSummarizedExperiment objects, phenotype tables, sample bigWigs, mean bigWigs, and file information tables are ready to use and freely available here. We also created the [recount](#) Bioconductor package which allows you to search and download the data for a specific study . By taking care of several preprocessing steps and combining many datasets into one easily-accessible website, we make finding and analyzing RNA-seq data considerably more straightforward.

Related publications

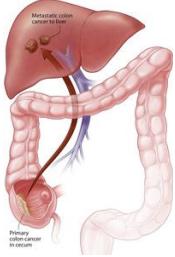
Collado-Torres L, Nellore A, Kammers K, Ellis SE, Taub MA, Hansen KD, Jaffe AE, Langmead B, Leek JT. [recount: A large-scale resource of analysis-ready RNA-seq expression data](#). *bioRxiv* **068478**.

The Datasets

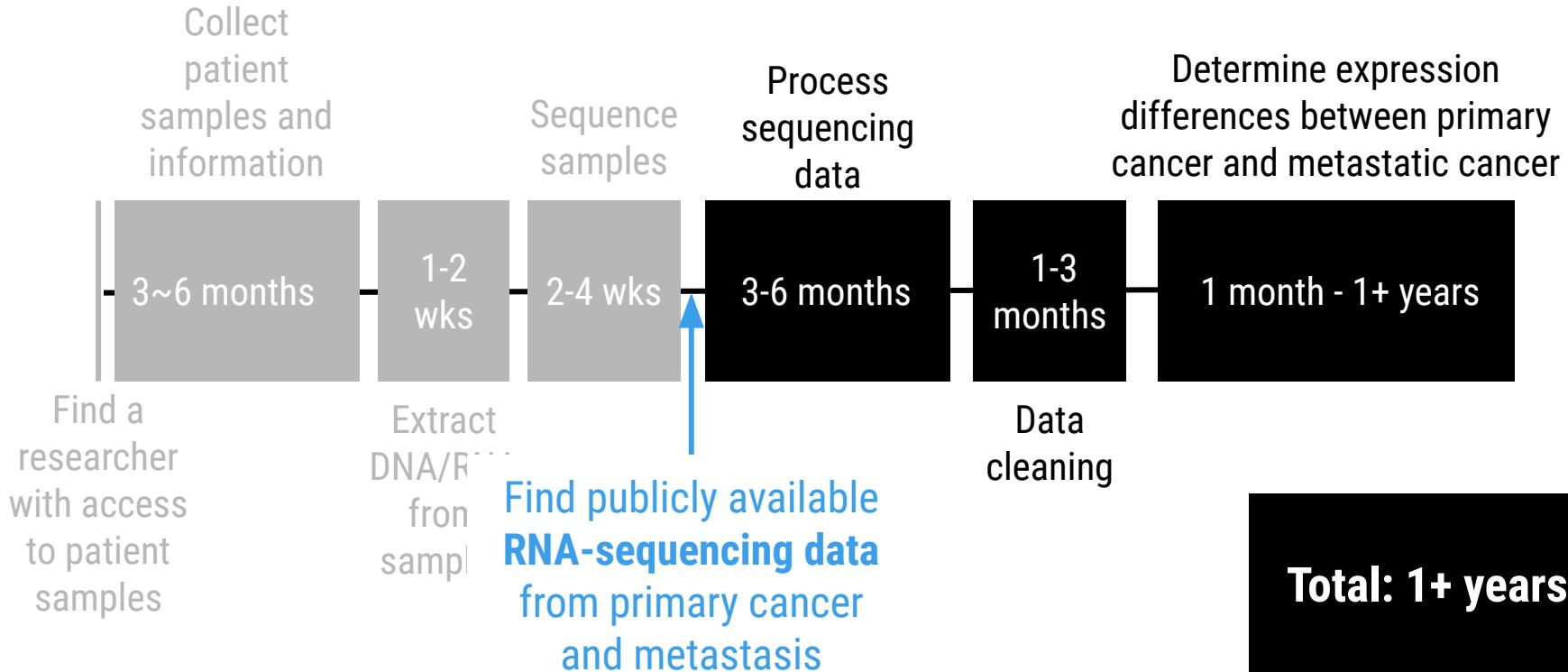
Show 10 ↑ entries

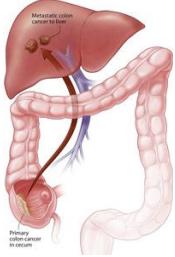
Search:

accession	number of samples	species	abstract	gene	exon	junctions	phenotype	info
All	All	All	All			All	All	
SRP025982	1720	human	We present primary results from the Sequencing Quality Control (SEQC) project, coordinated by the United States Food and Drug Administration. Examining Illumina HiSeq, Life Technologies SOLiD and Roche 454 platforms at multiple laboratory sites using reference RNA samples with built-in controls, we assess RNA sequencing (RNA-seq) performance for sequence discovery and differential expression profiling and compare it to microarray and quantitative PCR (qPCR) data using complementary metrics. At all sequencing depths, we discover unannotated exon-exon junctions, with >80% validated by qPCR. We find that	RSE counts	RSE counts	RSE jx_bed jx_cov counts	link	link

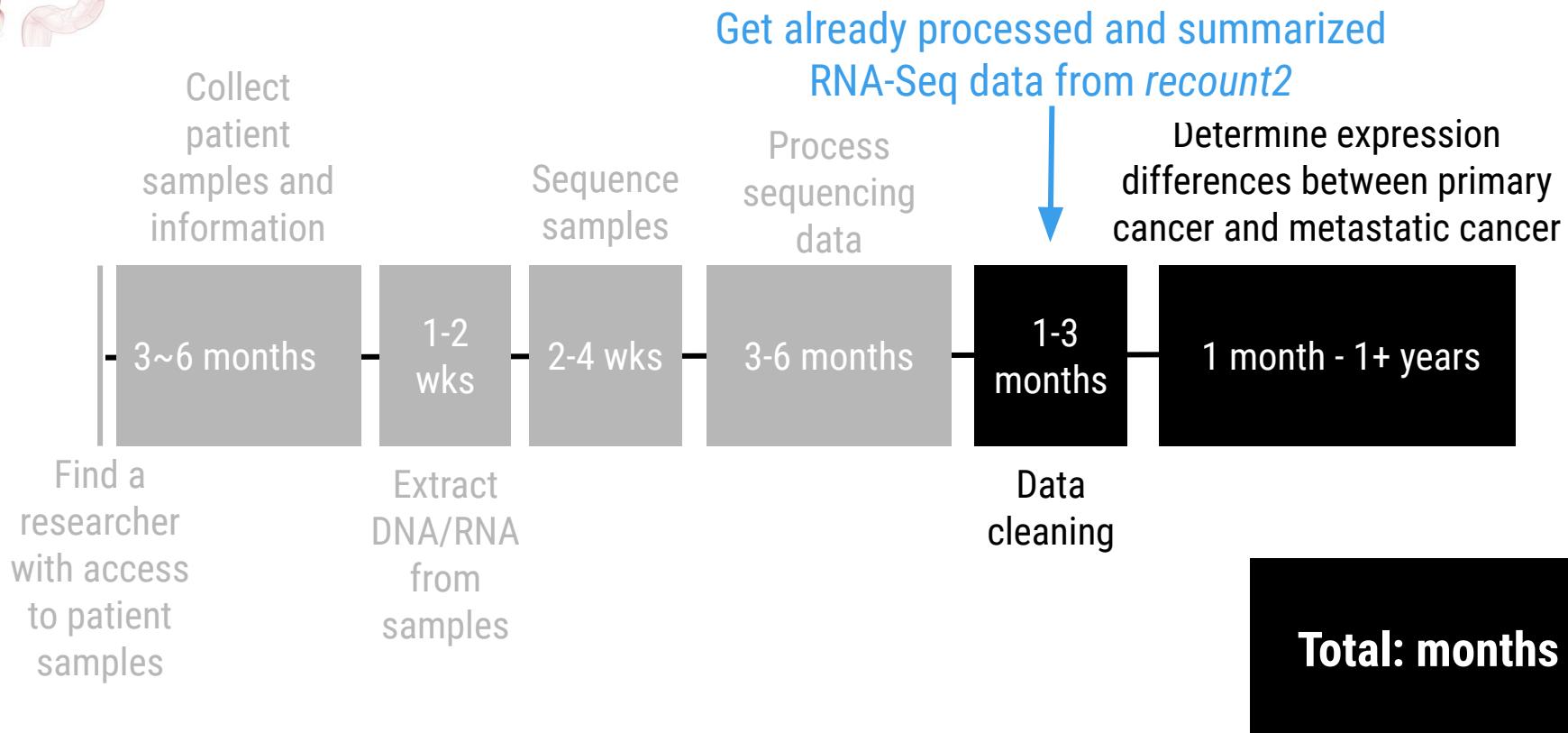


What makes primary cancer different than metastatic cancer?

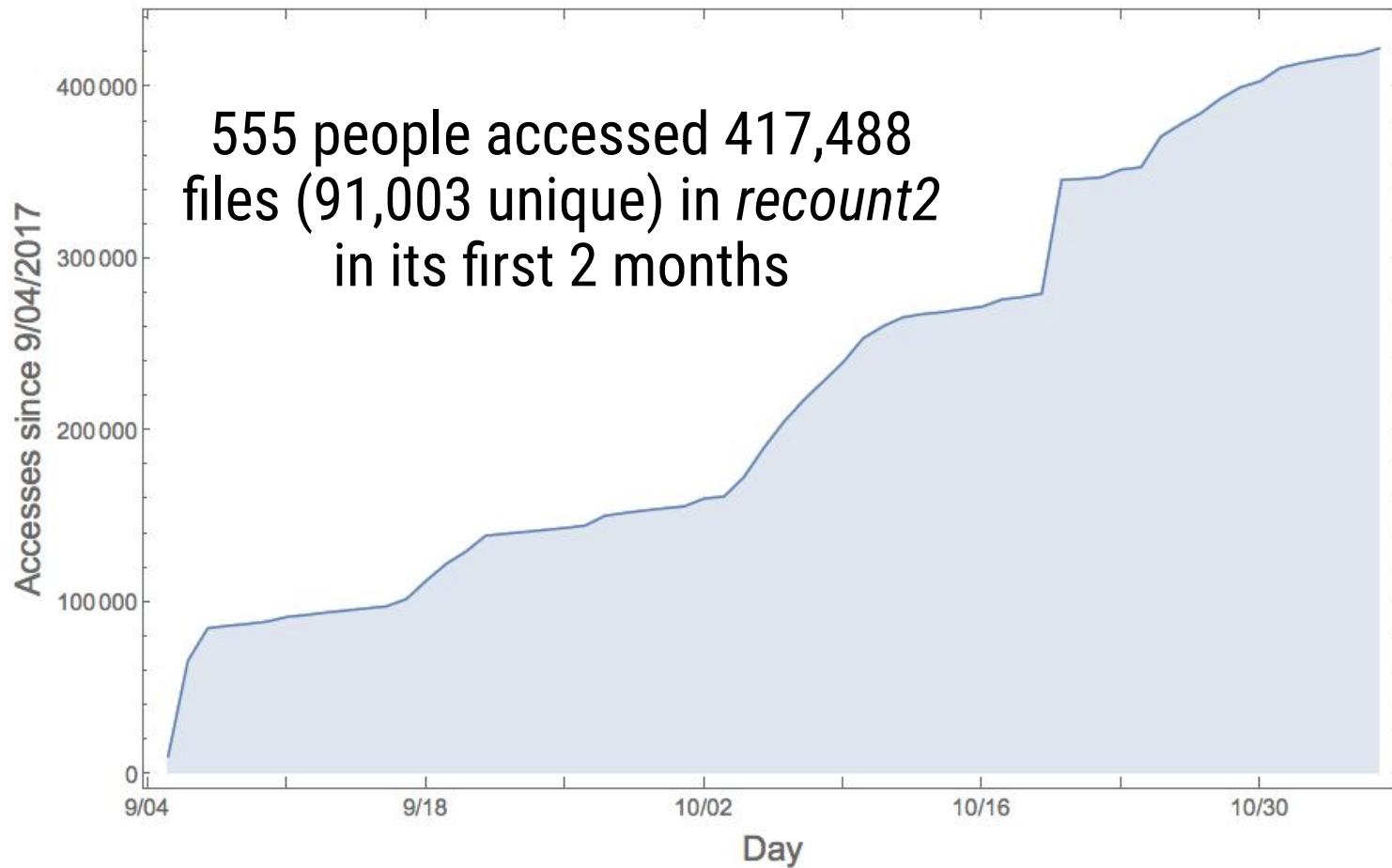




What makes primary cancer different than metastatic cancer?



recount accesses over two months of 2017



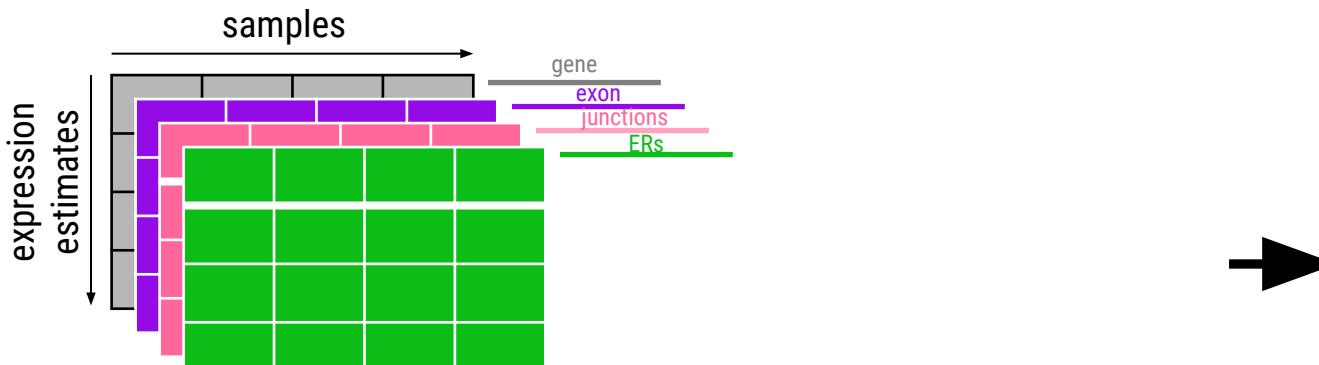
Project Goal: Take publicly available RNA-Seq data and make it available and easy-to-use

- Comprehensive
- Easy to Get
- Useful for future study



recount²

expression data for ~70,000 human samples

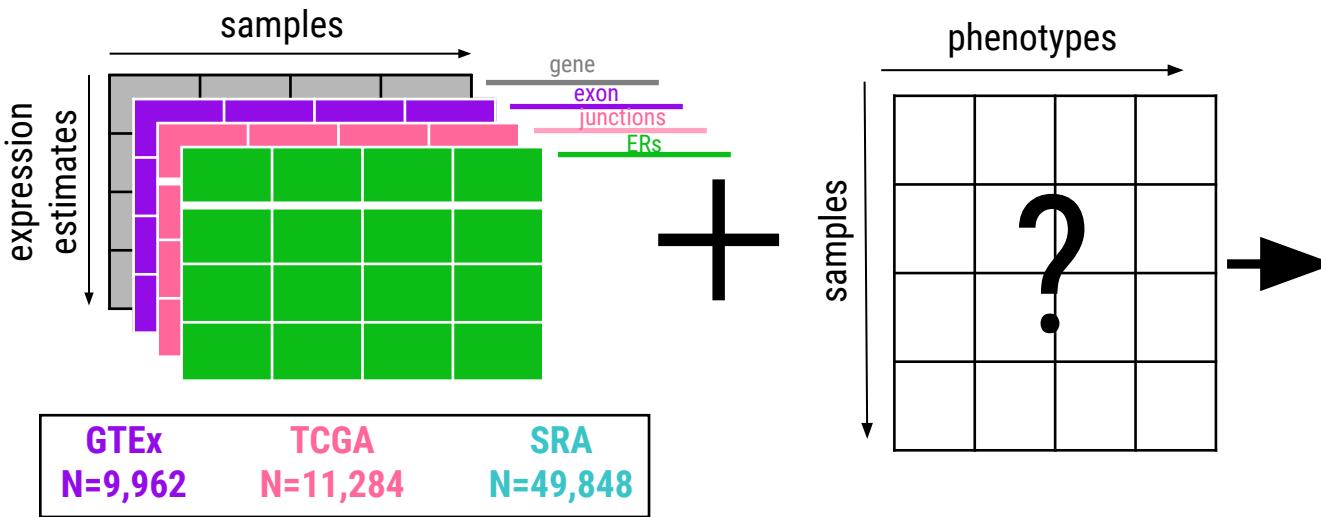


**Answer meaningful
questions about
human biology and
expression**

GTEX	TCGA	SRA
N=9,962	N=11,284	N=49,848

recount2

expression data for ~70,000 human samples



Answer meaningful questions about human biology and expression

SRA phenotype information is far from complete

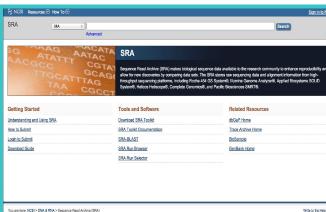
	Sex	Tissue	Race	Age
6620	female	liver	NA	NA
6621	female	liver	NA	NA
6622	female	liver	NA	NA
6623	female	liver	NA	NA
6624	female	liver	NA	NA
6625	male	liver	NA	NA
6626	male	liver	NA	NA
6627	male	liver	NA	NA
6628	male	liver	NA	NA
6629	male	liver	NA	NA
6630	male	liver	NA	NA
6631	NA	blood	NA	NA
6632	NA	blood	NA	NA
6633	NA	blood	NA	NA
6634	NA	blood	NA	NA
6635	NA	blood	NA	NA
6636	NA	blood	NA	NA



SRA

SRA phenotype information is far from complete

	Sex	Tissue	Race	Age
6620	female	liver	NA	NA
6621	female	liver	NA	NA
6622	female	liver	NA	NA
6623	female	liver	NA	NA
6624	female	liver	NA	NA
6625	male	liver	NA	NA
6626	male	liver	NA	NA
6627	male	liver	NA	NA
6628	male	liver	NA	NA
6629	male	liver	NA	NA
6630	male	liver	NA	NA
6631	NA	blood	NA	NA
6632	NA	blood	NA	NA
6633	NA	blood	NA	NA
6634	NA	blood	NA	NA
6635	NA	blood	NA	NA
6636	NA	blood	NA	NA



SRA

Even when information *is* provided, it's not always clear...

The screenshot shows the homepage of the Sequence Read Archive (SRA). At the top, there is a search bar and a navigation menu with links like "Home", "About", "Help", and "Logout". Below the header, there is a large orange box containing a sequence of DNA bases: AGTAGGAAATATAACGCCCTTGCAATTGGAGTAAAGCGCT. To the right of this sequence, there is a section titled "SRA" with a brief description of the archive's purpose. Below the sequence, there are three columns of links: "Getting Started", "Tools and Software", and "Related Resources". Under "Getting Started", links include "Introduction and Using SRA", "How to Submit", "Using SRA", and "Submit". Under "Tools and Software", links include "Overview SRA Tools", "SRA Toolkit Documentation", "SRA Toolkit", "SRA Toolkit", and "SRA Run Selector". Under "Related Resources", links include "SRA Home", "Sequence Read Archive", "SRA Home", "SRA Home", and "SRA Run Selector". At the bottom of the page, there is a footer with links to "National Center for Biotechnology Information (NCBI)", "National Library of Medicine (NLM)", "National Institutes of Health (NIH)", and "National Human Genome Research Institute (NHGRI)".

SRA

Category	Frequency
F	95
female	2036
Female	51
M	77
male	1240
Male	141
Total	3640

Even when information *is* provided, it's not always clear...



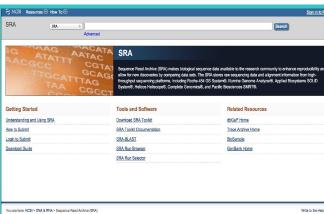
The screenshot shows the Sequence Read Archive (SRA) homepage. At the top, there's a search bar and a navigation menu with links like "Home", "About", "Help", and "Log In". Below the header, there's a large orange box containing a DNA sequence: "TG-AGAG-AATA-TA-GACGCC-TTGCATTTC-TAA-CGGCT". To the right of the sequence, there's a section titled "SRA" with a brief description of the archive. Below this are three main sections: "Getting Started", "Tools and Software", and "Related Resources". Under "Getting Started", there are links for "Introduction and Using SRA", "How to Submit", "Using SRA", and "Submit". Under "Tools and Software", there are links for "Overview SRA Tools", "SRA Toolkit Documentation", "SRA Toolkit", and "SRA Toolkit Select". Under "Related Resources", there are links for "SRA Home", "Sequence Read Archive", "SRA Toolkit", and "SRA Toolkit Select". At the bottom of the page, there's a footer with links for "Contact Us", "Help", "Log In", and "Logout".

SRA

Category	Frequency
F	95
female	2036
Female	51
M	77
male	1240
Male	141
Total	3640

"1 Male, 2 Female", "2 Male, 1 Female", "3 Female", "DK", "male and female" "Male (note:)", "missing", "mixed", "mixture", "N/A", "Not available", "not applicable", "not collected", "not determined", "pooled male and female", "U", "unknown", "Unknown"

Even when information *is* provided, it's not always clear...

A screenshot of the Sequence Read Archive (SRA) website. At the top, there's a search bar and a navigation menu with options like "Home", "About", "Help", and "Logout". Below the menu, there's a large orange box containing a DNA sequence: "TG-AGAG-AATA-TA GACGCC TTAGATTTC TAA CGCG". To the right of the sequence, there's a section titled "SRA" with a brief description: "Sequence Read Archive (SRA) collects biological sequence data available to the research community in archive repositories after the raw data has been quality checked. The SRA stores raw sequencing data and alignment information from SAGE, Illumina, Roche 454, ABI SOLiD, Ion Torrent, Solexa, and other sequencing platforms. The SRA is part of the National Biotechnology Information System (NBIS), which also includes the National Center for Biotechnology Information (NCBI), the National Library of Medicine (NLM), the National Institutes of Health (NIH), and the National Human Genome Research Institute (NHGRI).". Below this, there are three columns: "Getting Started" (with links to "Introduction and Using SRA", "How to Submit", "List of Submitters", and "Submitter Guide"), "Tools and Software" (with links to "Overview SRA Tools", "SRA Toolkit Documentation", "SRA Toolkit", "SRA Toolkit Basic", and "SRA Toolkit Select"), and "Related Resources" (with links to "dbSNP Home", "Sequence Read Archive", "National Center for Biotechnology Information", and "National Human Genome Research Institute"). At the bottom, there's a footer with links to "Help", "Logout", and "Search SRA".

SRA

Category	Frequency
F	95
female	2036
Female	51
M	77
male	1240
Male	141
Total	3640

"1 Male, 2 Female", "2 Male, 1 Female", "3 Female", "DK", "male and female" "Male (note:)", "missing", "mixed", "mixture", "N/A", "Not available", "not applicable", "not collected", "not determined", "pooled male and female", "U", "unknown", "Unknown"

# of NAs	# w/sex assigned
44,957	4,700

in-silico Phenotyping



Goal:
 to accurately
 predict critical
 phenotype
 information for
 all samples in
recount2



recount2
gene, exon, exon-exon junction and expressed region RNA-Seq data

GTEx

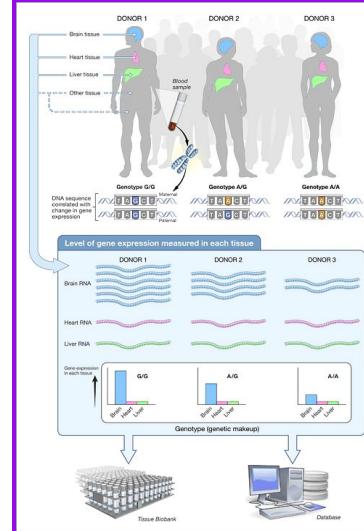
Genotype Tissue Expression Project
N=9,662

TCGA

The Cancer Genome Atlas
N=11,284

SRA

Sequence Read Archive
N=49,848



GTEx
 Genotype Tissue Expression Project
 N=9,662



TCGA
 The Cancer Genome Atlas
 N=11,284



SRA
 Sequence Read Archive
 N=49,848

Machine Learning: Making predictions from data

Data Set #1

Data Set #2

Machine Learning: Making predictions from data

Data Set #1

Data Set #2

Training
Data

Test
Data

Machine Learning: Making predictions from data

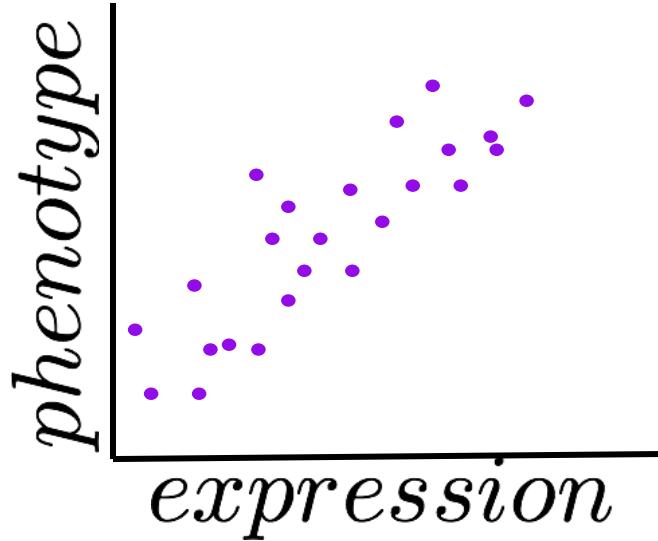
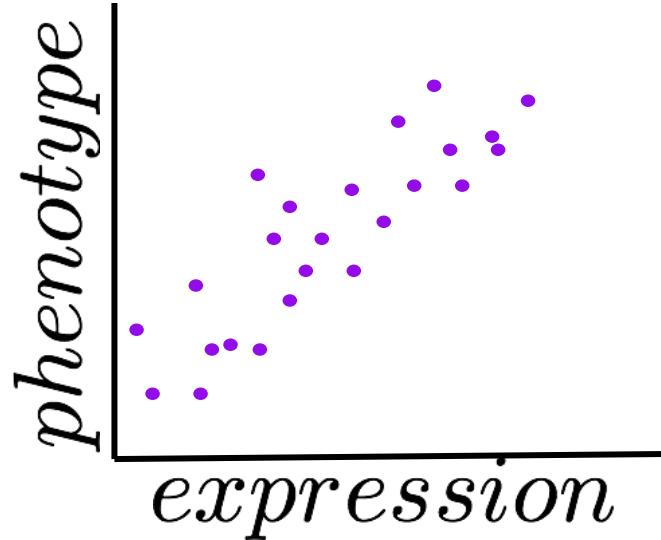
Data Set #1

Data Set #2

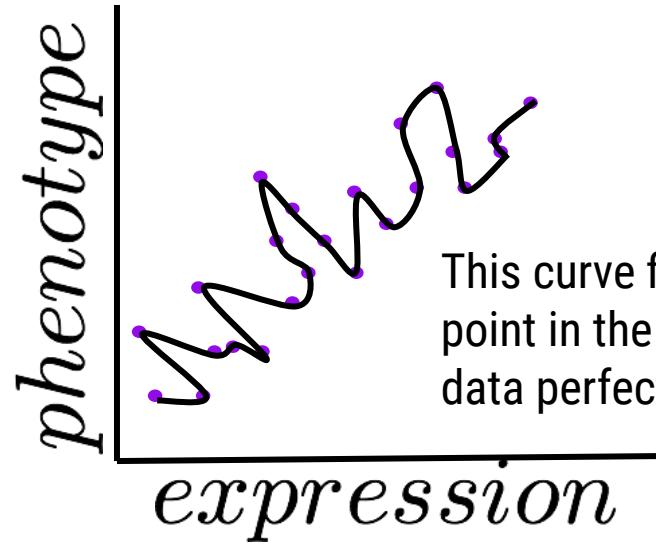
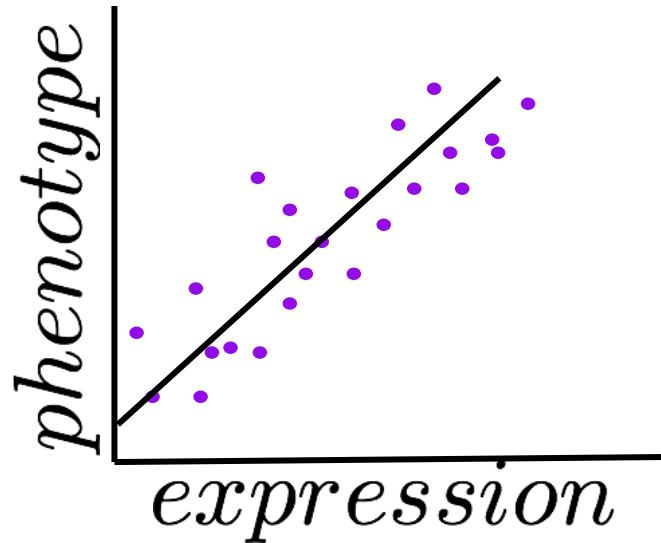
Training
Data

Data used to build
the predictor

We're interested in predicting phenotype from expression data...



There are a number of different curves you could fit through these data



Machine Learning: Making predictions from data

Data Set #1

Data Set #2

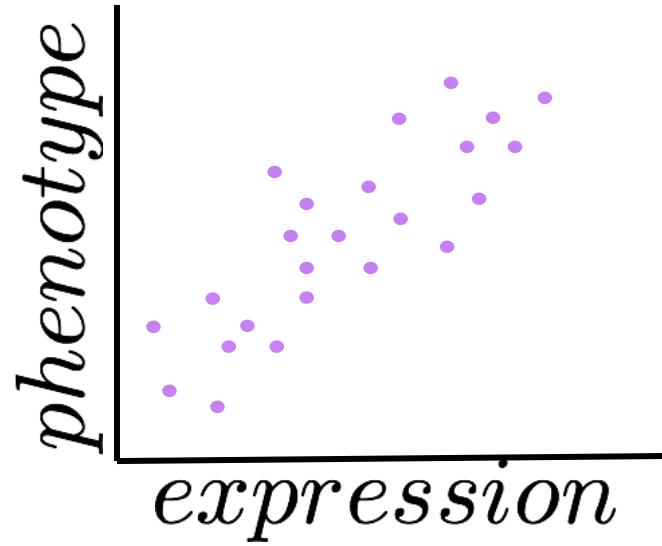
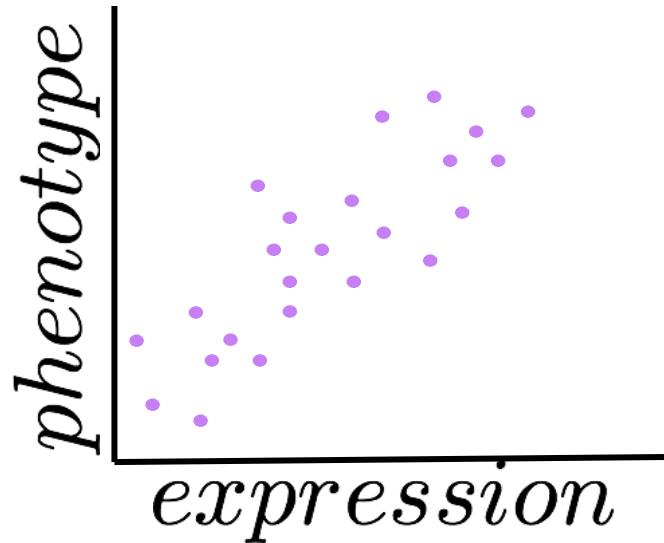
Training
Data

Test
Data

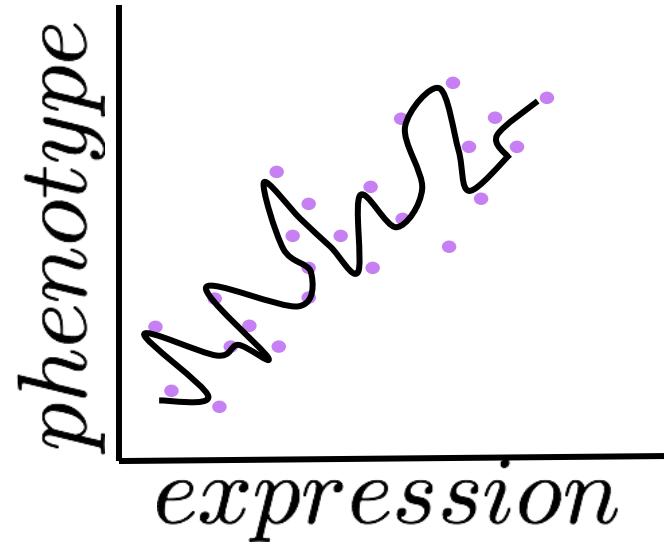
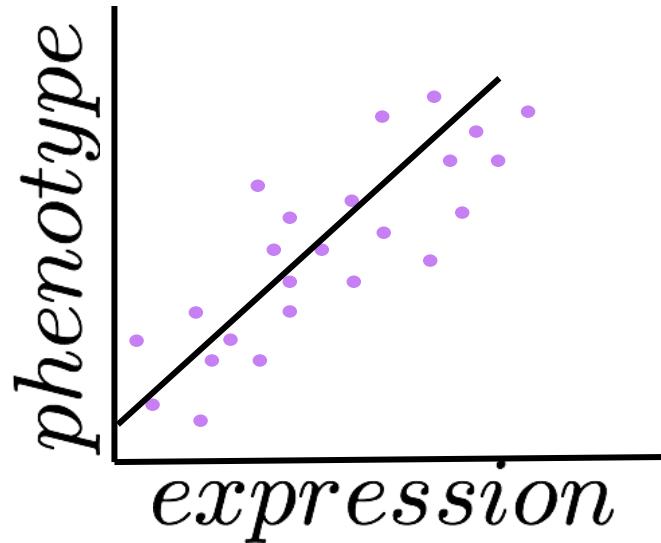
Data used to build
the predictor

Samples held
back from training

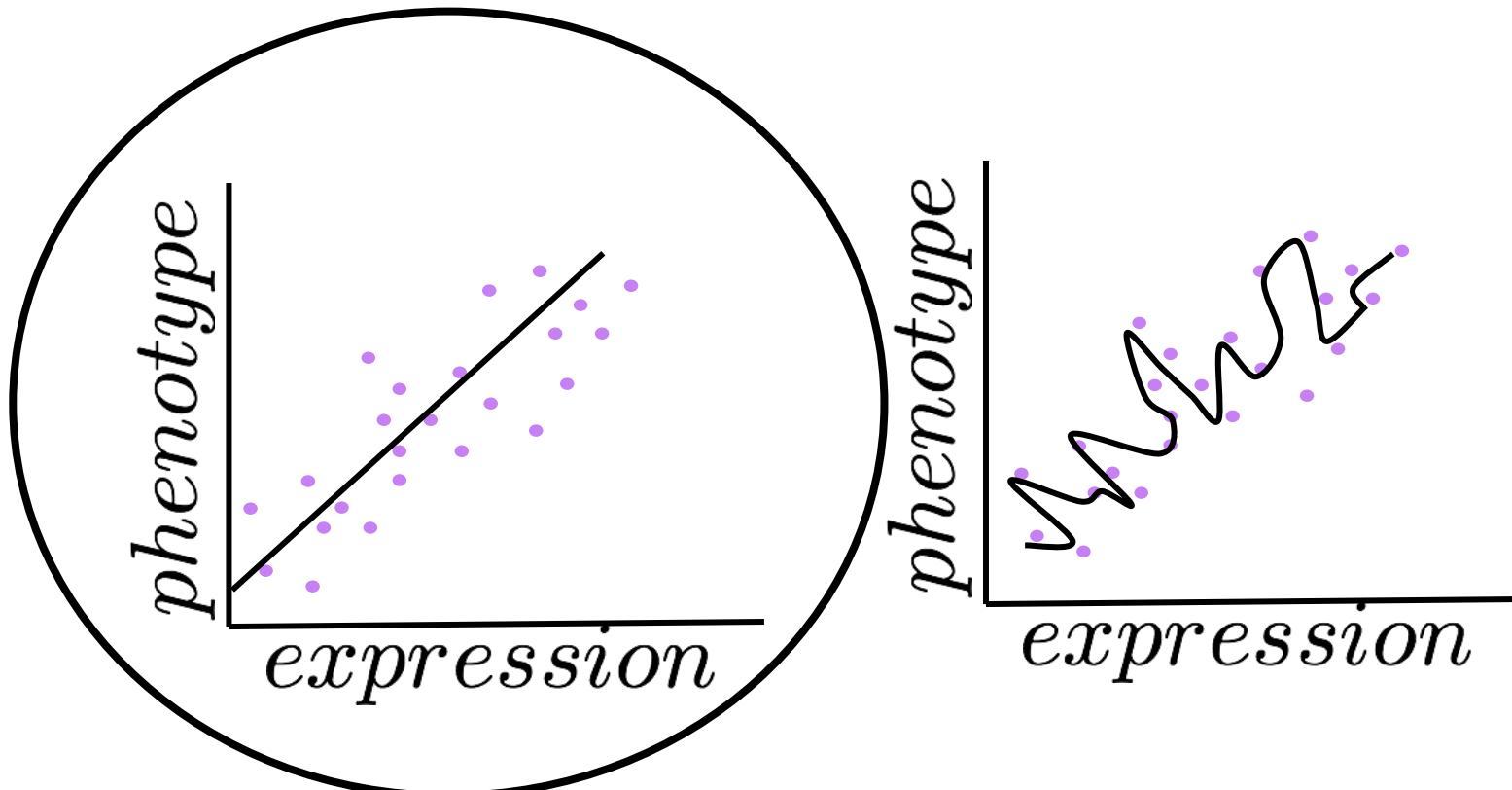
What if we tried to predict phenotype in the test data?



The curve no longer fits every point perfectly.



The curve no longer fits every point perfectly.



Machine Learning: Making predictions from data

Data Set #1

Data Set #2

Training
Data

Test
Data

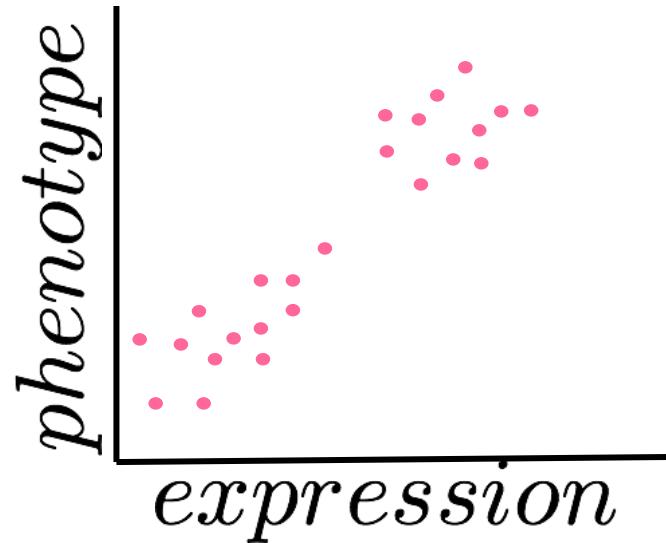
Validation
Data

Data used to build
the predictor

Samples held
back from training

Independent data set to
test predictor

We can now test prediction accuracy in an independent set of samples.



The line generated from the training data accurately predicts phenotype in the validation data



Prediction can be done using **linear regression**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

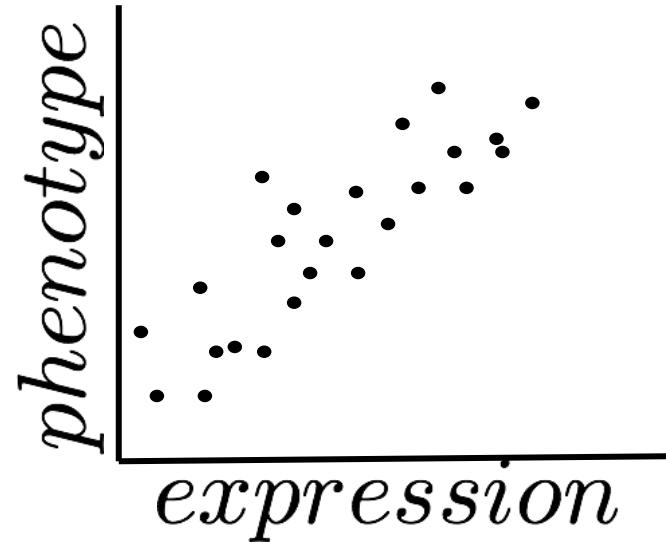
Prediction can be done using **linear regression**

$$Y = \beta_0 + \beta_1 X + \epsilon$$

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$

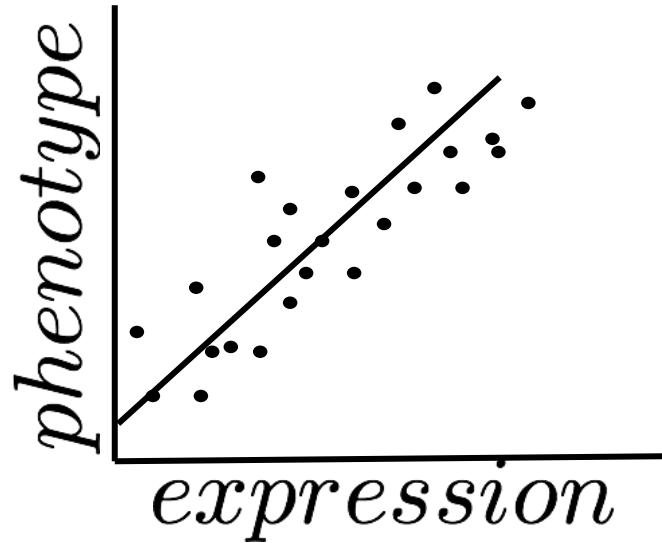
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



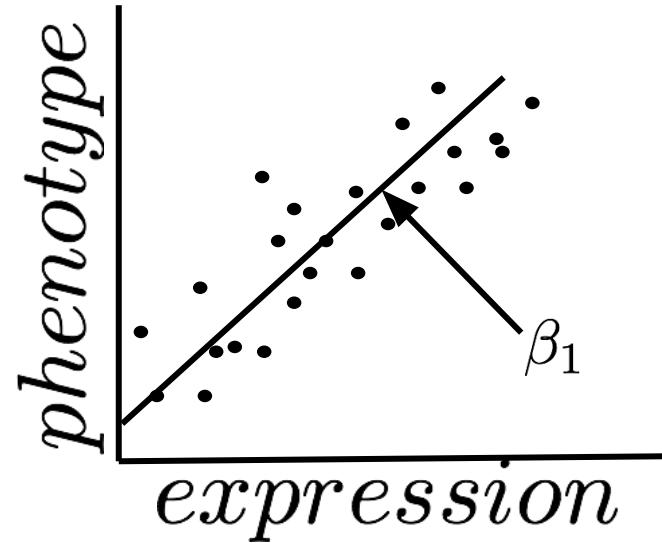
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



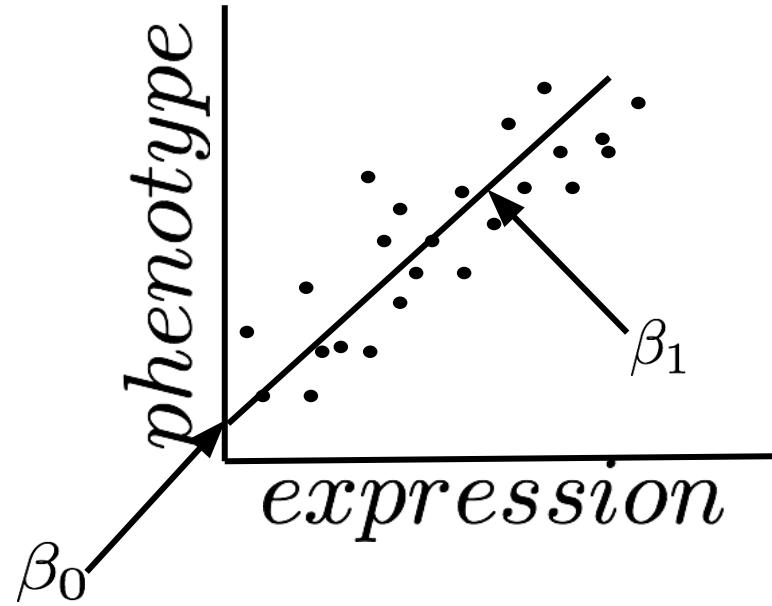
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



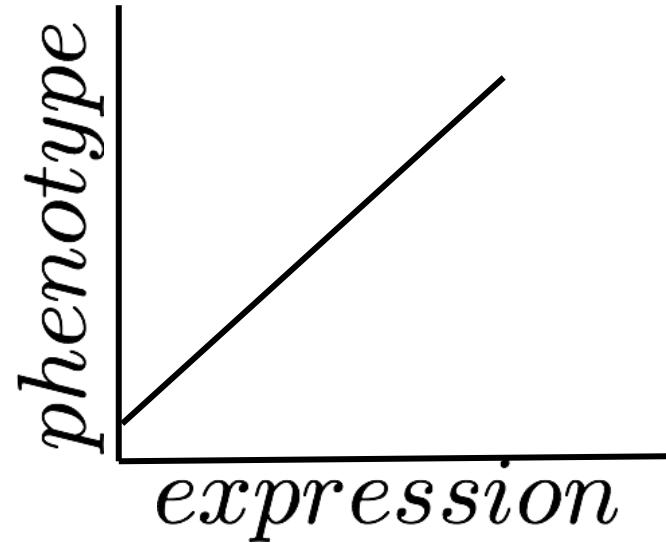
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



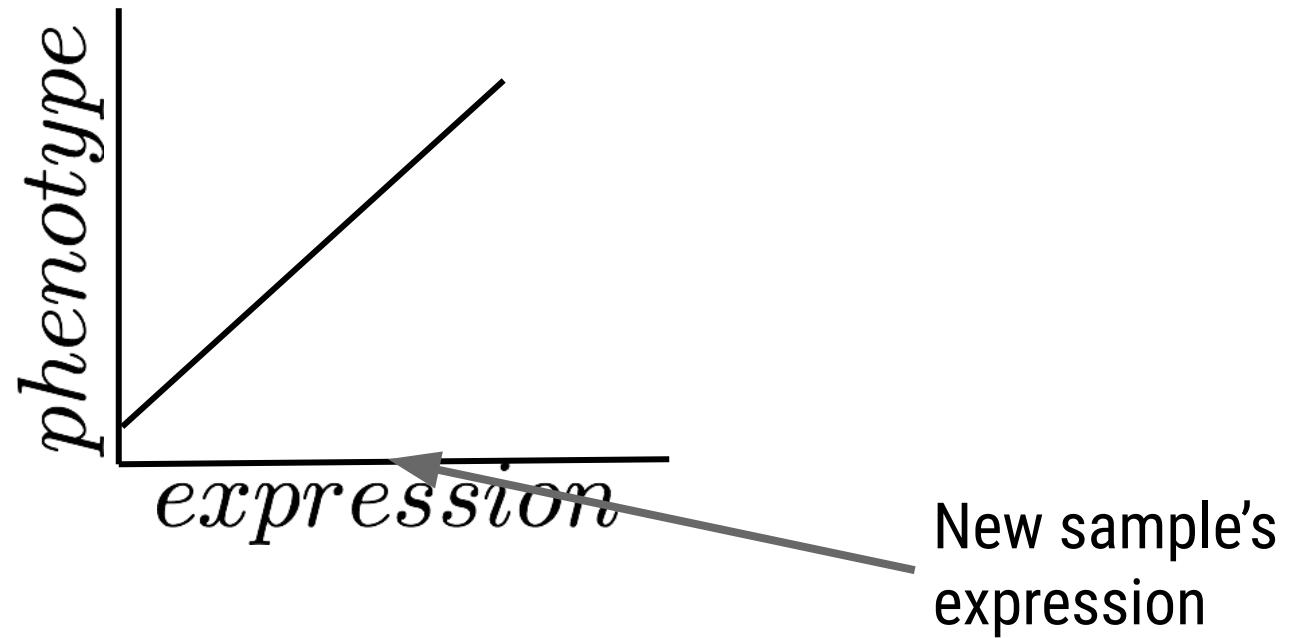
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



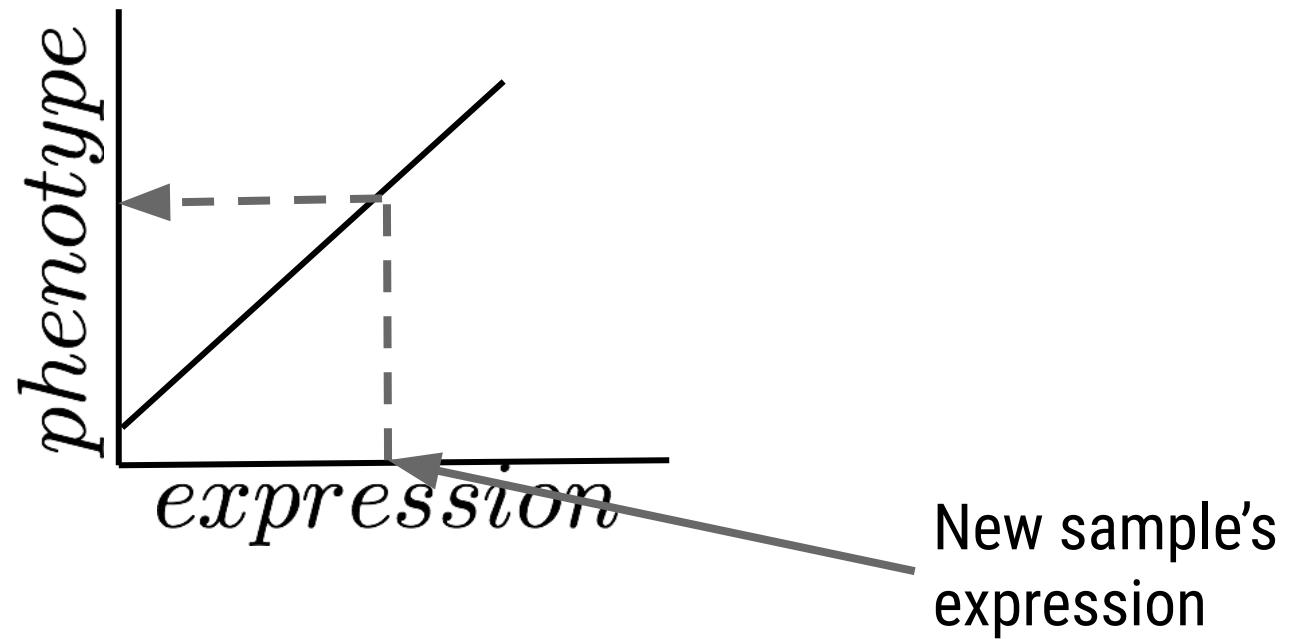
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



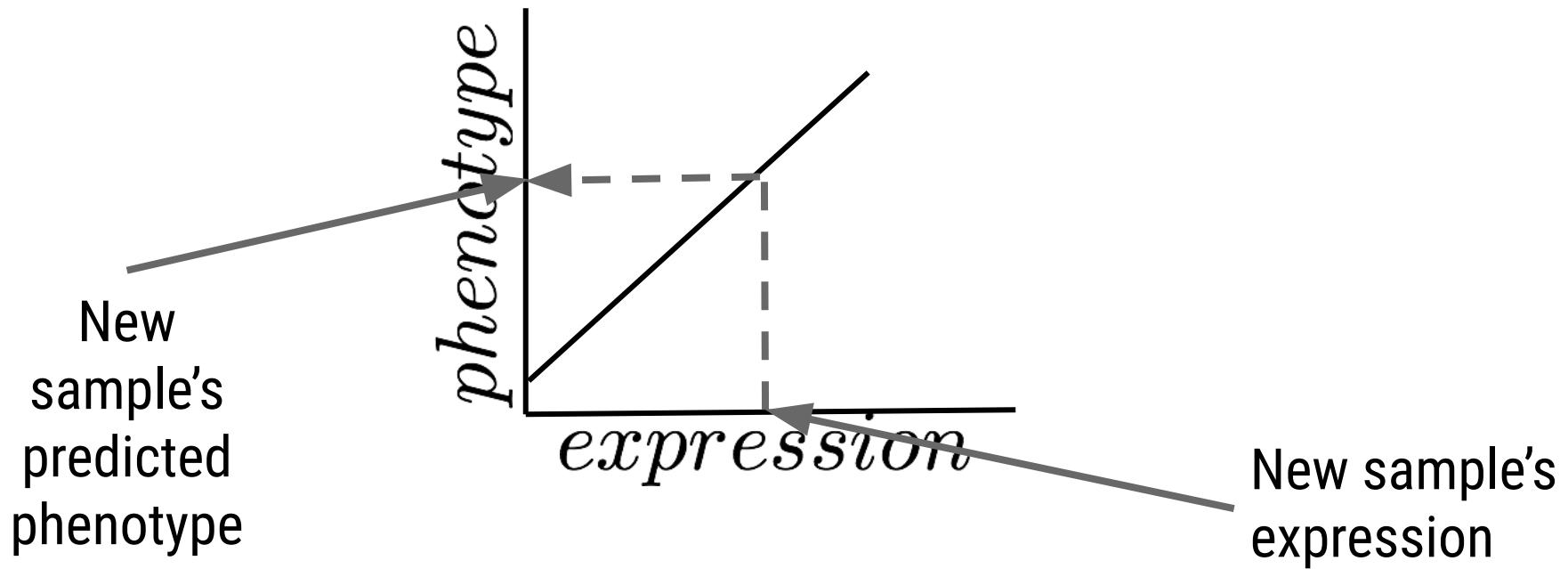
Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



Prediction is done using **linear regression**

$$\text{phenotype} = \beta_0 + \beta_1(\text{expression}) + \epsilon$$



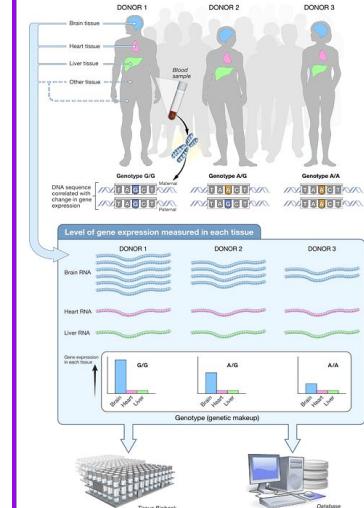
Goal:
 to accurately
 predict critical
 phenotype
 information for
 all samples in
recount2



GTEX
Genotype Tissue Expression Project
 N=9,662

TCGA
The Cancer Genome Atlas
 N=11,284

SRA
Sequence Read Archive
 N=49,848



NATIONAL CANCER INSTITUTE THE CANCER GENOME ATLAS

TCGA BY THE NUMBERS

- TCGA produced over **2.5 PETABITES** of data
- TCGA data describes **33 TUMOR TYPES**, including **10 CANCERS** based on paired tumor and normal tissue samples collected from **212,000 PATIENTS** using **7 DIFFERENT DATA TYPES**

TCGA RESULTS & FINDINGS

- NEED TO KNOW THE BASIS OF CANCER
- TCGA SURVEYS
- IDENTIFYING TARGETS
- THE TEAM

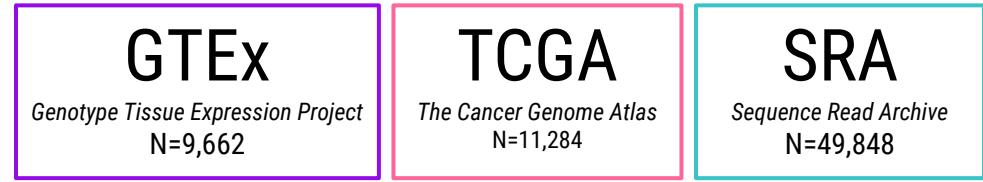
WHAT'S NEXT?

The Genome Data Commons (GDC) will host TCGA data along with other NCIT-generated data and facilitate access from anywhere. This will allow researchers to more easily expand their analyses to allow researchers to ask more relevant questions with TCGA data.

www.cancer.gov/ccg

Goal :

to accurately
predict critical
phenotype
information for
all samples in
recount2

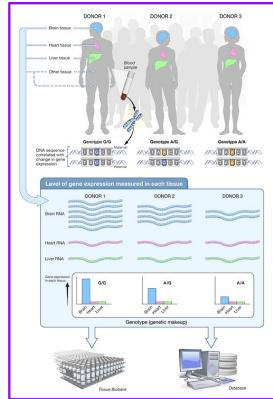


build and optimize
phenotype
predictor

Training
Data

Missingness limited in GTEx phenotype data

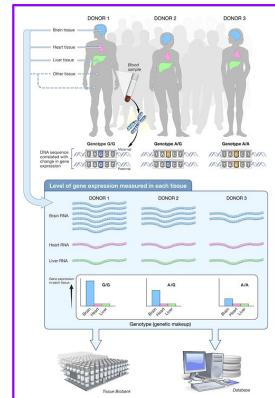
	Sex	Tissue	Race	Age
1	male	Lung	White	59
2	male	Brain	White	27
3	female	Heart	Black or African American	23
4	male	Brain	White	51
5	male	Skin	White	27
6	male	Lung	White	68
7	female	Brain	White	61
8	female	Adipose Tissue	White	42
9	male	Brain	White	40
10	female	Uterus	White	33
11	female	Nerve	White	60
12	male	Muscle	White	54
13	female	Ovary	White	31
14	male	Blood	White	53
15	female	Brain	White	56
16	male	Muscle	White	44



GTEx

Missingness limited in GTEx phenotype data

	Sex	Tissue	Race	Age
1	male	Lung	White	59
2	male	Brain	White	27
3	female	Heart	Black or African American	23
4	male	Brain	White	51
5	male	Skin	White	27
6	male	Lung	White	68
7	female	Brain	White	61
8	female	Adipose Tissue	White	42
9	male	Brain	White	40
10	female	Uterus	White	33
11	female	Nerve	White	60
12	male	Muscle	White	54
13	female	Ovary	White	31
14	male	Blood	White	53
15	female	Brain	White	56
16	male	Muscle	White	44

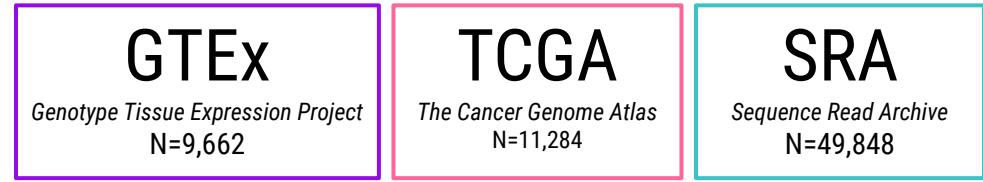


GTEx

Category	Frequency
female	3,626
male	6,036
NA	0

Goal :

to accurately
predict critical
phenotype
information for
all samples in
recount2

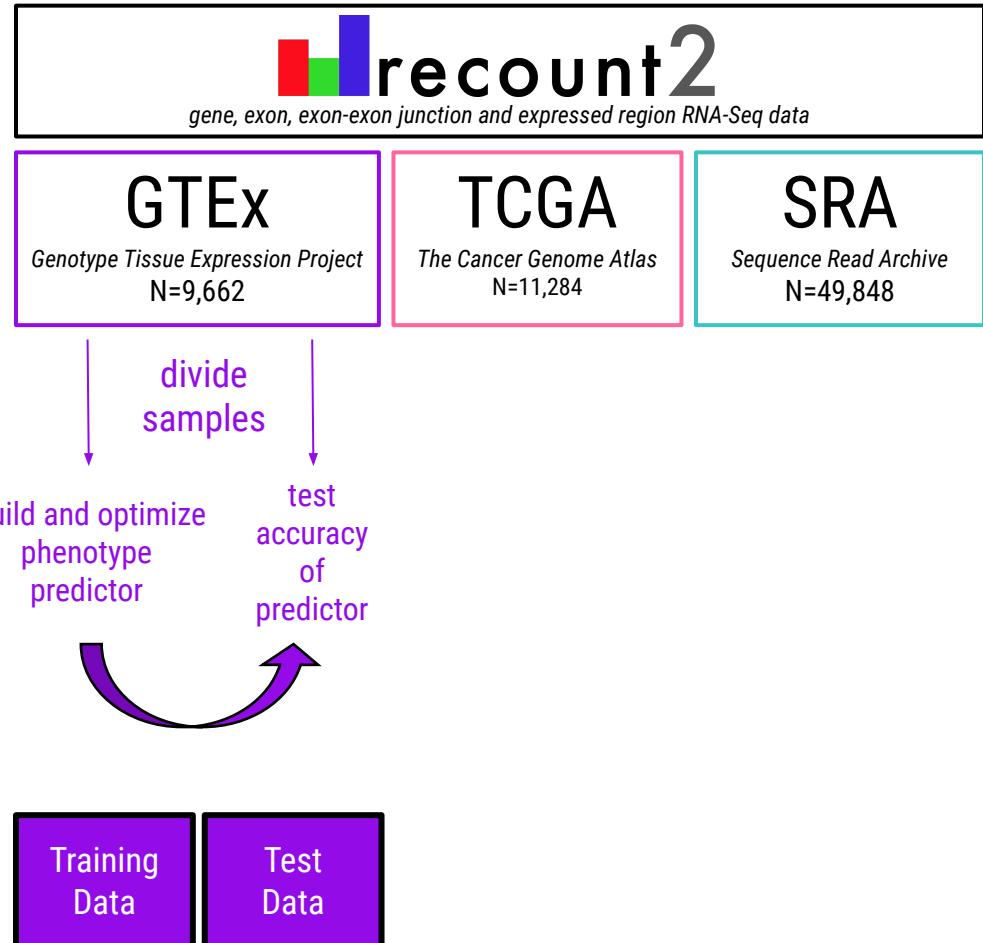


build and optimize
phenotype
predictor



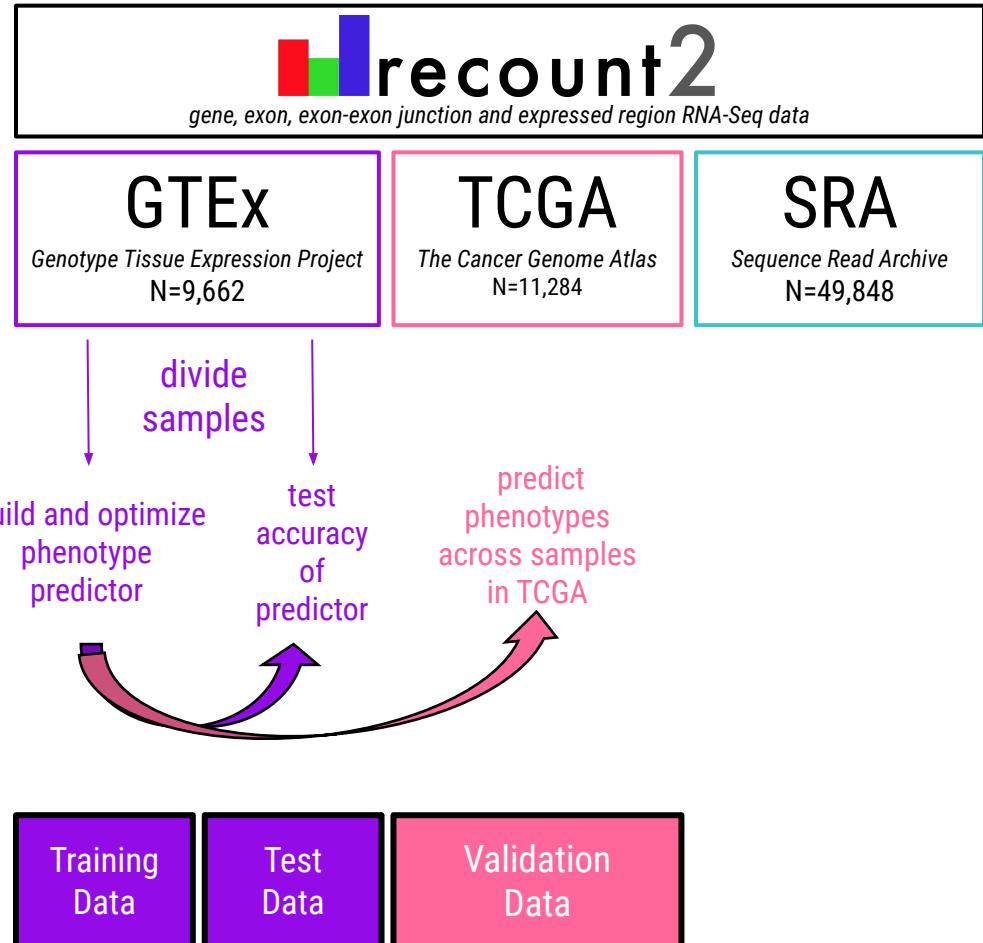
Goal :

to accurately predict critical phenotype information for all samples in *recount2*



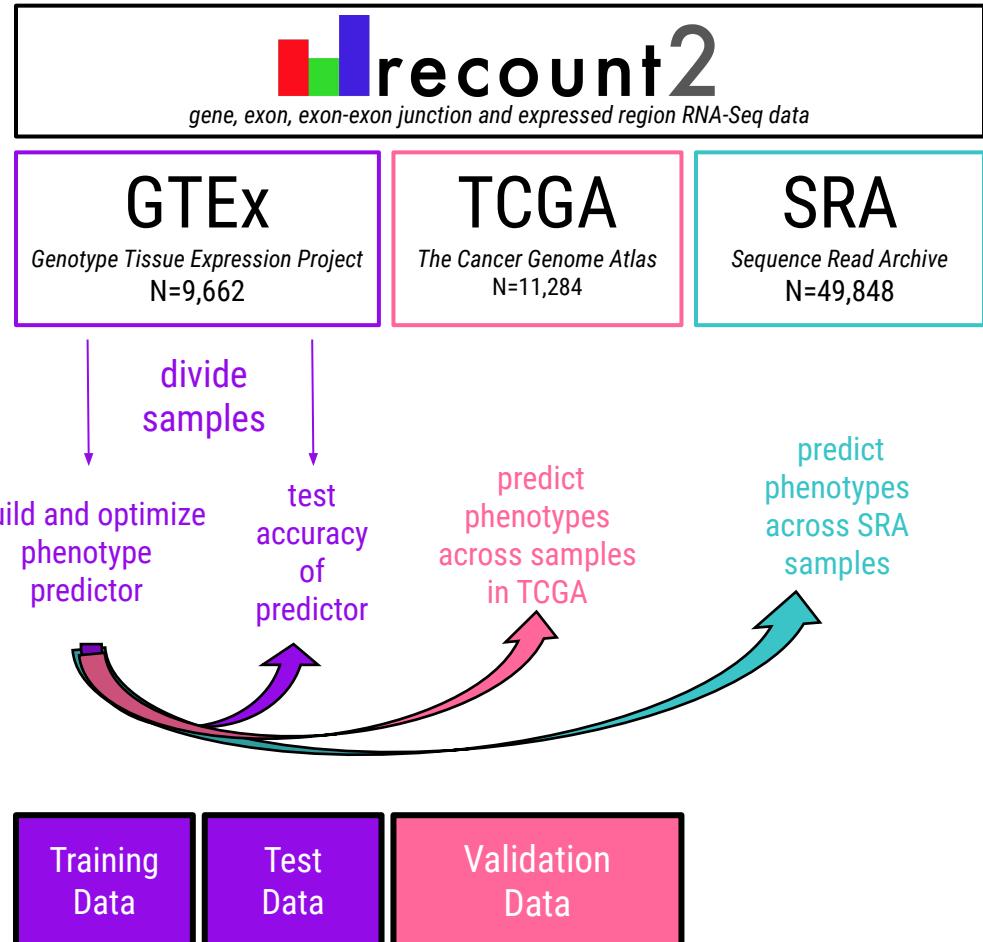
Goal :

to accurately predict critical phenotype information for all samples in *recount2*



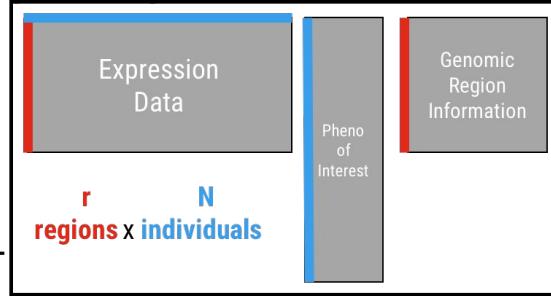
Goal :

to accurately predict critical phenotype information for all samples in *recount2*



phenopredict

Input Data

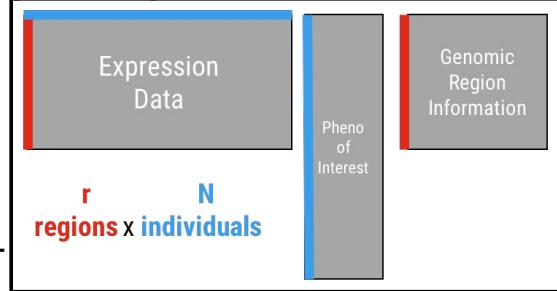


functions

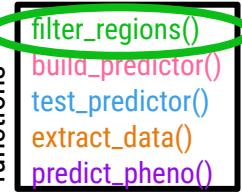


phenopredict

Input Data



functions



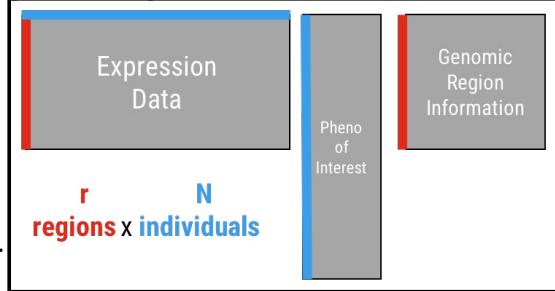
filter_regions()

Identify regions with differential expression for each level



phenopredict

Input Data



functions



filter_regions()

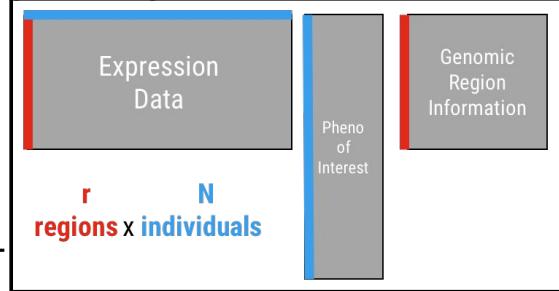
Identify regions with differential expression for each level



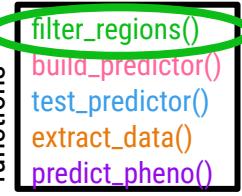
Set of discriminatory regions

phenopredict

Input Data

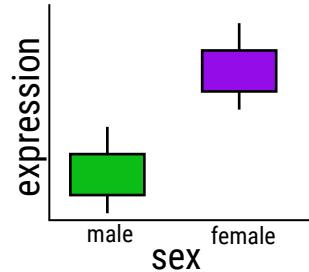


functions



filter_regions()

Identify regions with differential expression for each level



Set of discriminatory regions

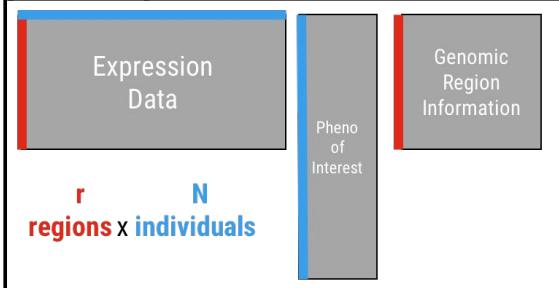
$$E[\mathbf{E}_r | \mathbf{P}_l] = \alpha_{0r} + \mathbf{P}_l \alpha_r$$

\mathbf{E}_r = expression at region r

$$\mathbf{P}_l = \begin{cases} 1, & \text{if level of phenotype} \\ 0, & \text{otherwise} \end{cases}$$

phenopredict

Input Data

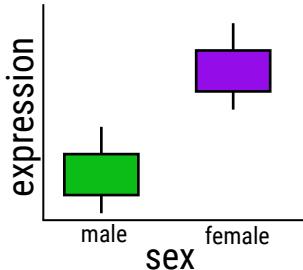


functions:

- filter_regions()
- build_predictor()
- test_predictor()
- extract_data()
- predict_pheno()

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

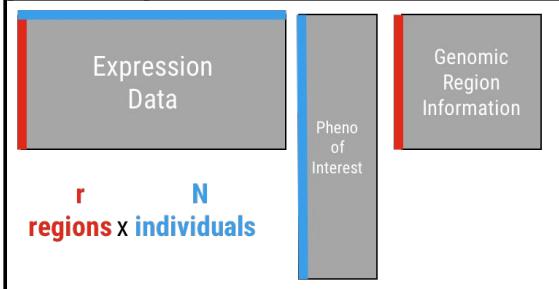
β_r phenotype (P_i)

male	female

Set of discriminatory regions

phenopredict

Input Data



functions:

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_i)

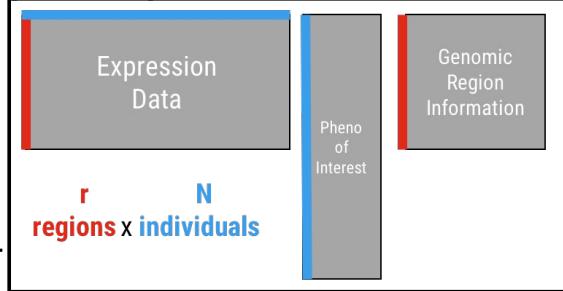
filtered regions (r)	male	female
1		
2		
3		
4		
5		

Set of discriminatory regions

Relationship between expression and phenotype

phenopredict

Input Data

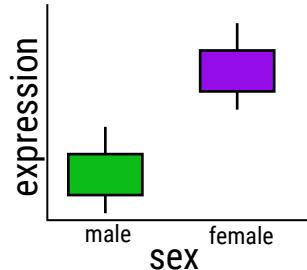


functions:

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

$$\beta_r \quad \text{phenotype } (P_i)$$

filtered regions (\mathcal{S})	male	female

Set of discriminatory regions

Relationship between expression and phenotype

$$E[\{\mathbf{E}_r\}_{r \in \mathcal{S}} | \mathbf{P}] = \boldsymbol{\beta}_r \mathbf{P}$$

$$\beta_r > 0 \forall \beta_r$$

$\mathbf{P} = l \times N$ matrix

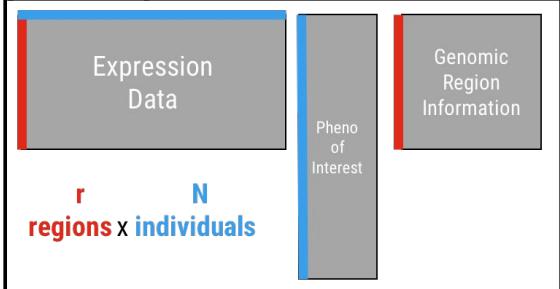
l = number of levels in the phenotype

N = number of samples

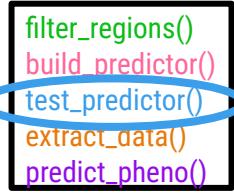
\mathbf{E} = expression for every region (r) in the selected set of regions (\mathcal{S})

phenopredict

Input Data



functions



filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_{rj} phenotype (P_j)

male	female

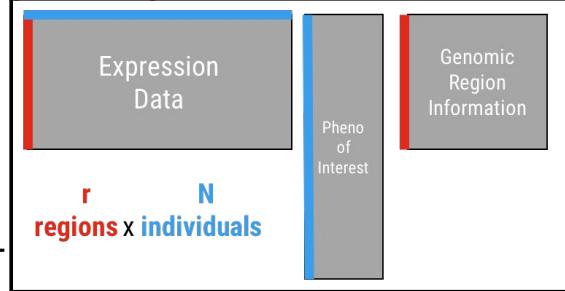
Set of discriminatory regions

Relationship between expression and phenotype

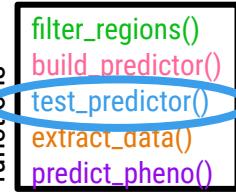
Likelihood of phenotype for each individual

phenopredict

Input Data



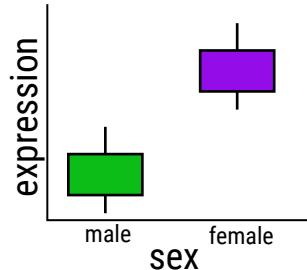
functions



test_predictor()

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_l)

male	female

Set of discriminatory regions

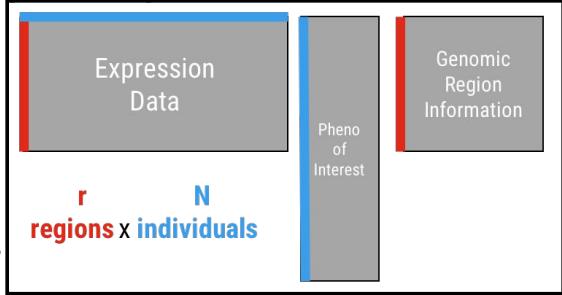
Relationship between expression and phenotype

$$E[\{\mathbf{E}_r\}_{r \in \mathcal{S}} | \hat{\beta}_r] = \hat{\beta}_r \gamma$$

Likelihood of phenotype for each individual

$\hat{\beta}_r$ = estimates of mean expression for each level (l) in phenotype (P)

phenopredict

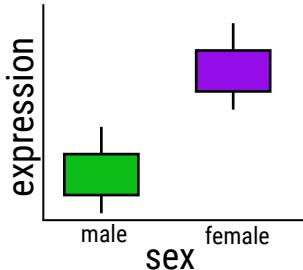


functions

```
filter_regions()  
build_predictor()  
test_predictor()  
extract_data()  
predict_pheno()
```

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

expression ~ phenotype

βr phenotype (P_1)

Set of discriminatory regions

Relationship between expression and phenotype

male	0.99	0.02	0.04	0.98
female	0.01	0.98	0.96	0.02

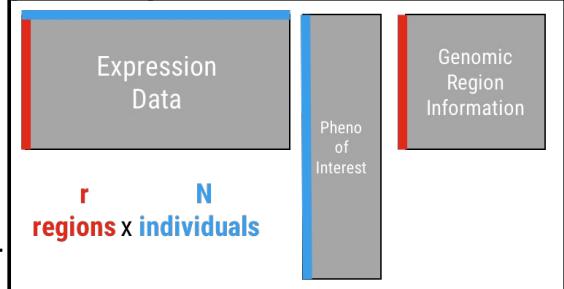
$$E[\{\mathbf{E}_r\}_{r \in \mathcal{S}} \mid \hat{\beta}_r] = \hat{\beta}_r \gamma$$

$\hat{\beta}_r$ = estimates of mean expression for each level (l) in phenotype (\mathbf{P})

$\hat{\gamma}$ = $l \times N$ matrix containing the likelihood of belonging to each level in the phenotype

phenopredict

Input Data



functions

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

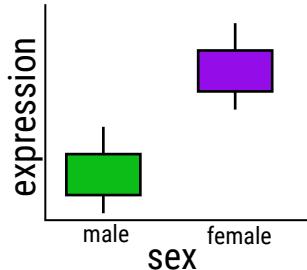
`test_predictor()`

Predict phenotype and assess accuracy in training set data

`predictions`

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_i)

	male	female
1		
2		
3		
4		

filtered regions (t)

Relationship between expression and phenotype

Set of discriminatory regions

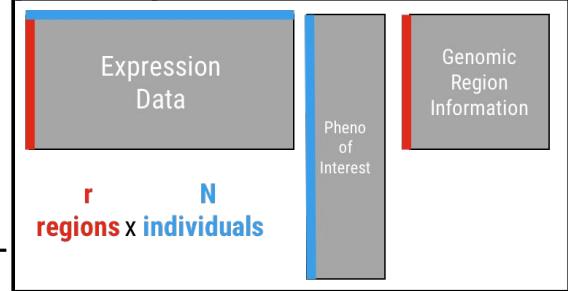
male	0.99	0.02	0.04	0.98
female	0.01	0.98	0.96	0.02

$\max_l(\hat{\gamma})$

Likelihood of phenotype for each individual

phenopredict

Input Data



functions

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

`test_predictor()`

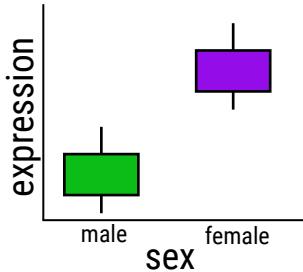
Predict phenotype and assess accuracy in training set data

`predictions`

Assign phenotype to most likely category

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_i)

	male	female
1		
2		
3		
4		

filtered regions (t)

Relationship between expression and phenotype

Set of discriminatory regions

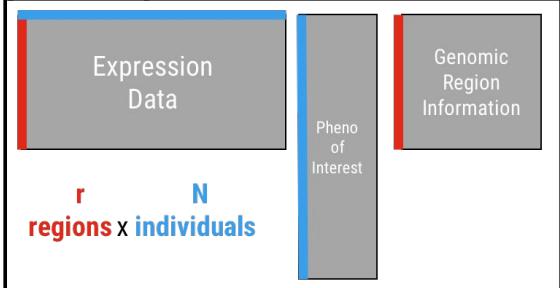
male	0.99	0.02	0.04	0.98
female	0.01	0.98	0.96	0.02

$\max_l(\hat{\gamma})$

Likelihood of phenotype for each individual

phenopredict

Input Data



functions

filter_regions()
build_predictor()
test_predictor()
extract_data()
predict_pheno()

test_predictor()

Predict phenotype and assess accuracy in training set data

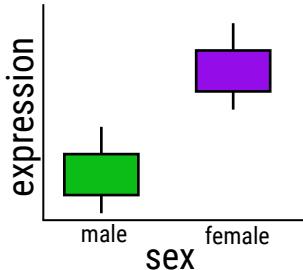
predictions

male				

Assign phenotype to most likely category

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_{rt} phenotype (P_i)

	male	female
male		
female		
male		
female		

filtered regions (t)

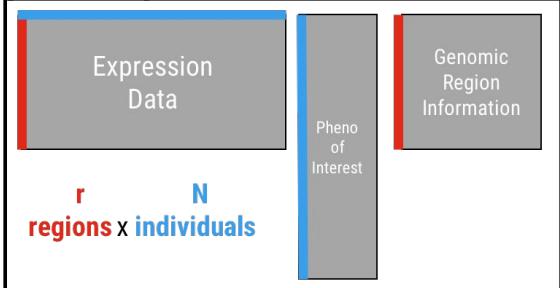
Set of discriminatory regions

male	0.99	0.02	0.04	0.98
female	0.01	0.98	0.96	0.02

Relationship between expression and phenotype

phenopredict

Input Data



functions

filter_regions()
build_predictor()
test_predictor()
extract_data()
predict_pheno()

test_predictor()

Predict phenotype and assess accuracy in training set data

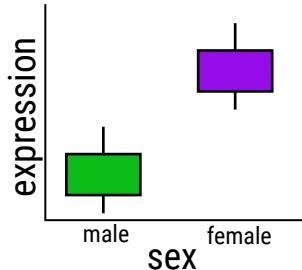
predictions

male
female
female
male

Assign phenotype to most likely category

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

expression ~ phenotype

β_r phenotype (P_i)

	male	female
male		
female		
male		
female		

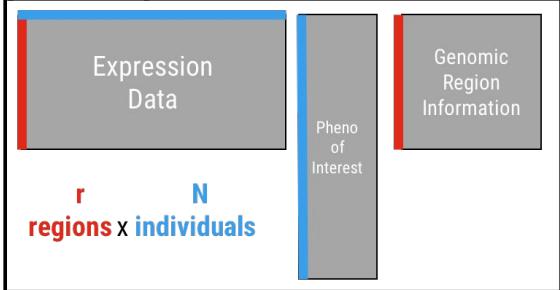
filtered regions (r)

Set of discriminatory regions

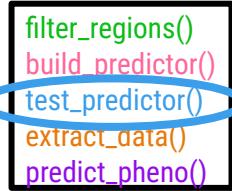
male	0.99	0.02	0.04	0.98
female	0.01	0.98	0.96	0.02

phenopredict

Input Data



functions



test_predictor()

Predict phenotype and assess accuracy in training set data

predictions

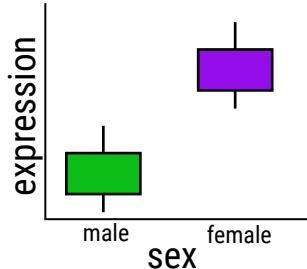
male
female
female
male

reported

male
female
female
male

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$expression \sim phenotype$

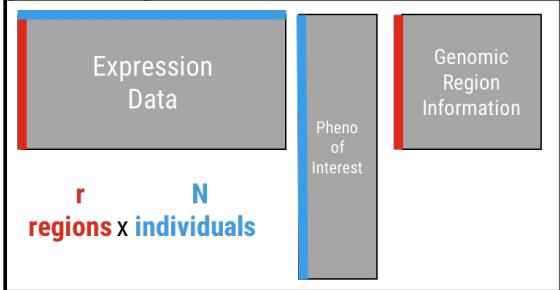
β_{rj} phenotype (P_j)

	male	female
1		
2		
3		
4		

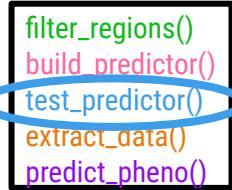
filtered regions (t)

phenopredict

Input Data



functions



test_predictor()

Predict phenotype and assess accuracy in training set data

predictions

male
female
female
male

reported

male
female
female
male

Prediction accuracy: 100%

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

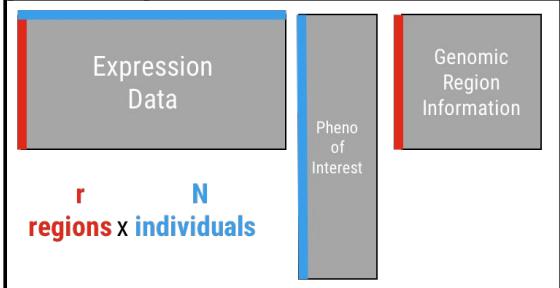
β_{rt} phenotype (P_t)

male female

filtered regions (t)

phenopredict

Input Data



functions

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

`test_predictor()`

Predict phenotype and assess accuracy in training set data

`predictions`

male
female
female
male

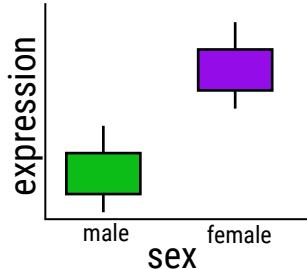
`reported`

male
female
female
male

Prediction accuracy: 100%

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

$\beta_{rt} \text{ phenotype } (P_t)$

male	female

filtered regions (t)

`extract_data()`

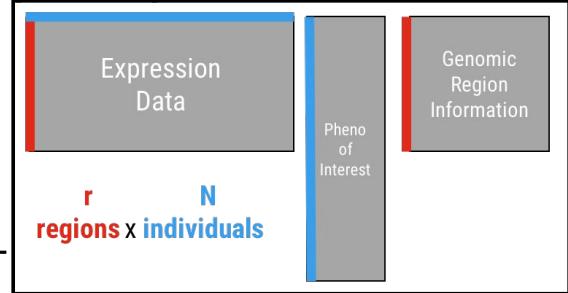
Extract expression information at regions identified by `filter_regions()` in a new data set

new data set samples

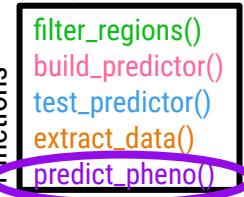
expression @
filtered regions

phenopredict

Input Data



functions



test_predictor()

Predict phenotype and assess accuracy in training set data

predictions

male
female
female
male

reported

male
female
female
male

Prediction accuracy: 100%

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_{rj} phenotype (P_j)

male	female

filtered regions (j)

predict_pheno()

Predict phenotypes across samples in this new data set

apply coefficient estimates to the extracted data

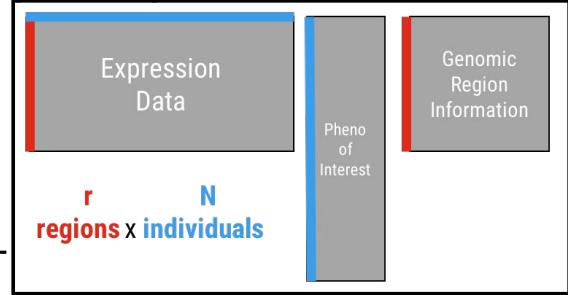
extract_data()

Extract expression information at regions identified by filter_regions() in a new data set

new data set samples	expression @	filtered regions

phenopredict

Input Data



functions

`filter_regions()`
`build_predictor()`
`test_predictor()`
`extract_data()`
`predict_pheno()`

`test_predictor()`

Predict phenotype and assess accuracy in training set data

`predictions`

male
female
female
male

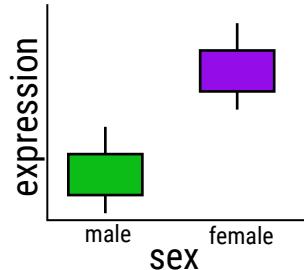
`reported`

male
female
female
male

Prediction accuracy: 100%

`filter_regions()`

Identify regions with differential expression for each level



`build_predictor()`

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_{rj} phenotype (P_j)

	male	female
1		
2		
3		
4		

filtered regions (r)

`predict_pheno()`

Predict phenotypes across samples in this new data set

`extract_data()`

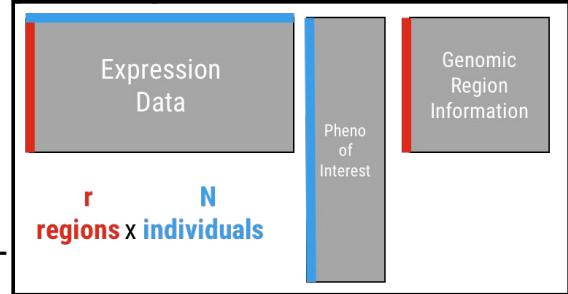
Extract expression information at regions identified by `filter_regions()` in a new data set

new data set samples	expression @ filtered regions
1	10, 15, 20, 25
2	25, 30, 35, 40
3	10, 15, 20, 25
4	25, 30, 35, 40

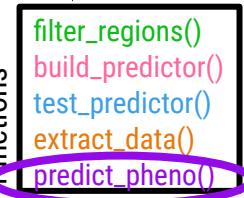
apply coefficient estimates to the extracted data

phenopredict

Input Data



functions



test_predictor()

Predict phenotype and assess accuracy in training set data

predictions

male
female
female
male

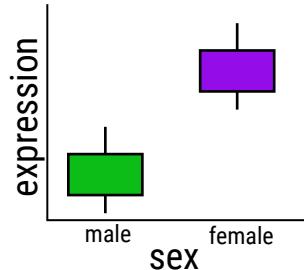
reported

male
female
female
male

Prediction accuracy: 100%

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_i)

male female

filtered regions (i)	male	female
1		
2		
3		
4		

extract_data()

Extract expression information at regions identified by filter_regions() in a new data set



predict_pheno()

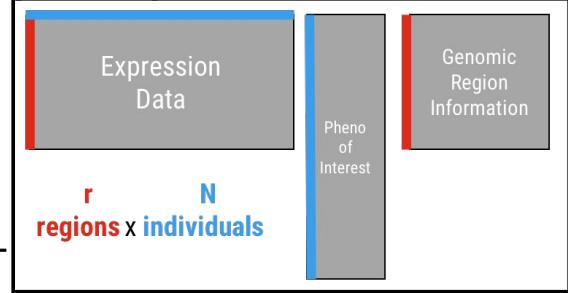
Predict phenotypes across samples in this new data set

apply coefficient estimates to the extracted data

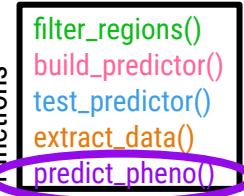
$$E[\{\mathbf{E}_r^*\}_{r \in \mathcal{S}} | \hat{\beta}_r] = \hat{\beta}_r \gamma^*$$

phenopredict

Input Data



functions



test_predictor()

Predict phenotype and assess accuracy in training set data

predictions

male
female
female
male

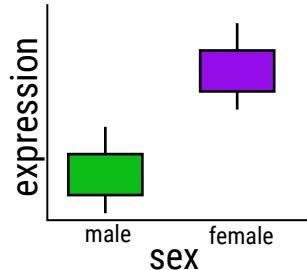
reported

male
female
female
male

Prediction accuracy: 100%

filter_regions()

Identify regions with differential expression for each level



build_predictor()

Extract coefficient estimates across regions

$\text{expression} \sim \text{phenotype}$

β_r phenotype (P_i)

male	female

filtered regions (t)

extract_data()

Extract expression information at regions identified by filter_regions() in a new data set



predict_pheno()

Predict phenotypes across samples in this new data set

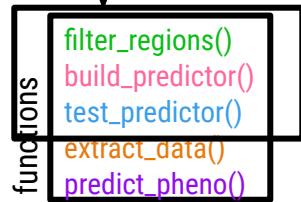
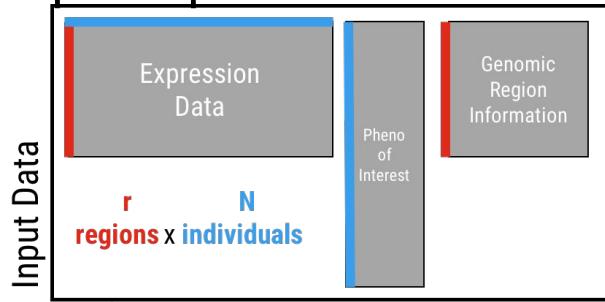
predictions in new data set

male
male
male
female

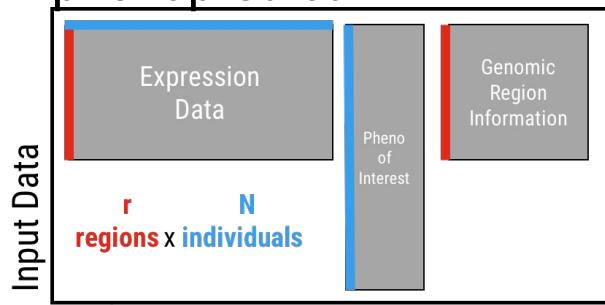
apply coefficient estimates to the extracted data

$$E[\{\mathbf{E}_r^*\}_{r \in \mathcal{S}} | \hat{\beta}_r] = \hat{\beta}_r \gamma^*$$

phenopredict

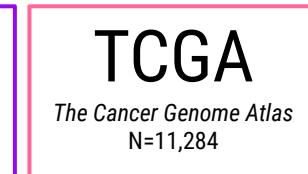
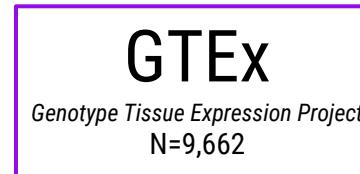
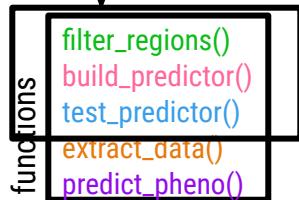


phenopredict

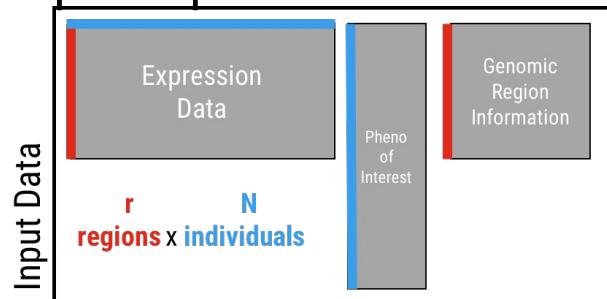


Accuracy

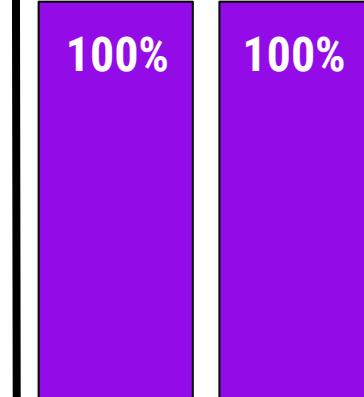
100%



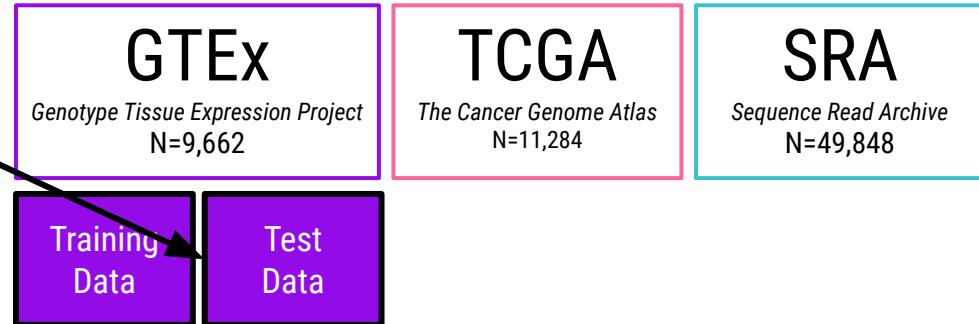
phenopredict



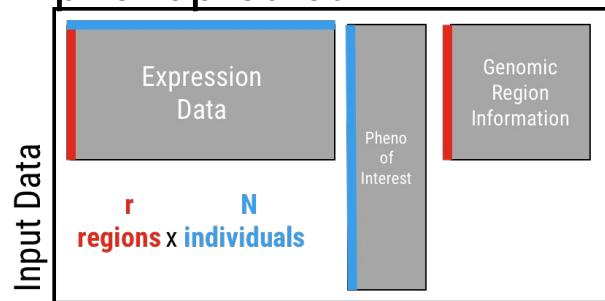
Accuracy



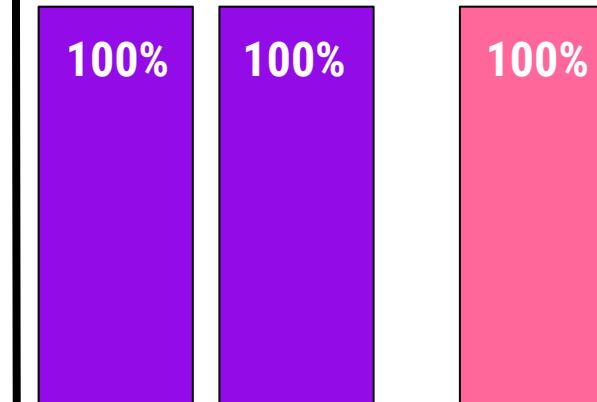
functions



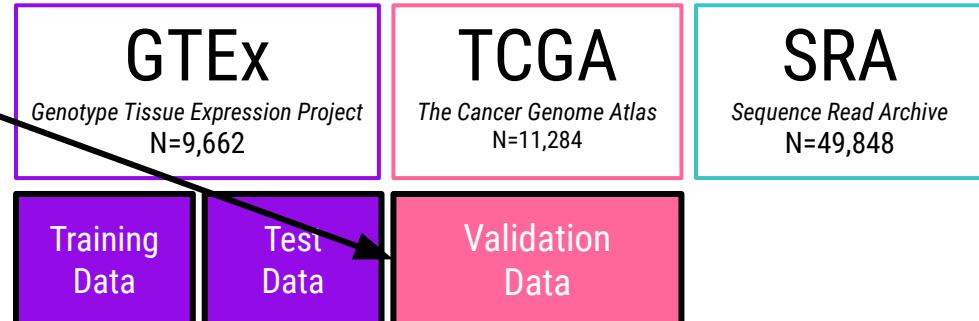
phenopredict



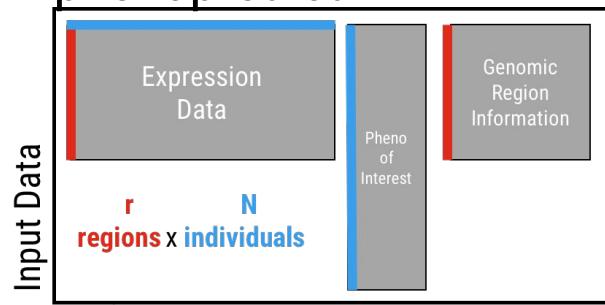
Accuracy



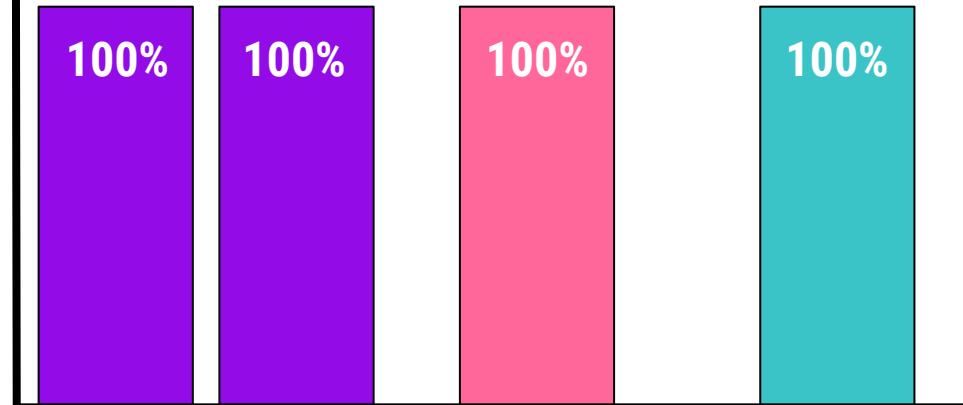
functions



phenopredict

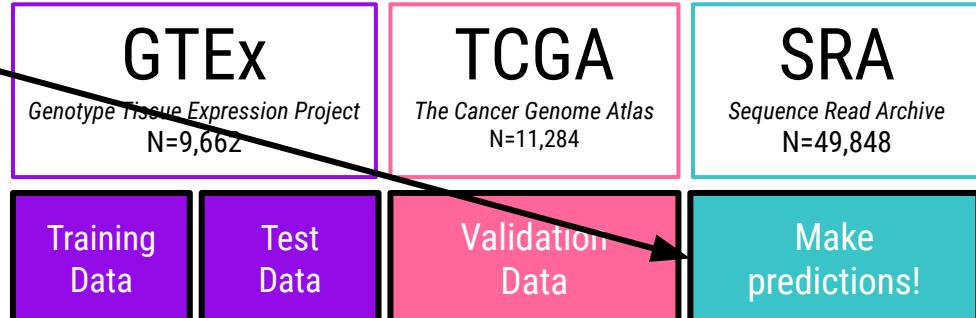


Accuracy

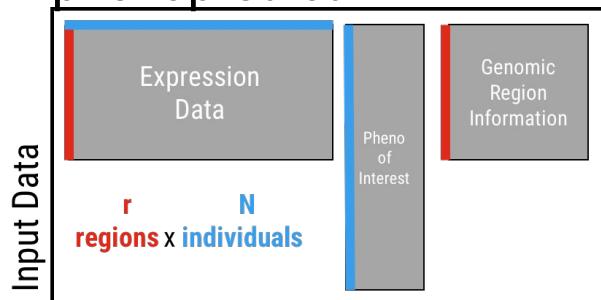


functions

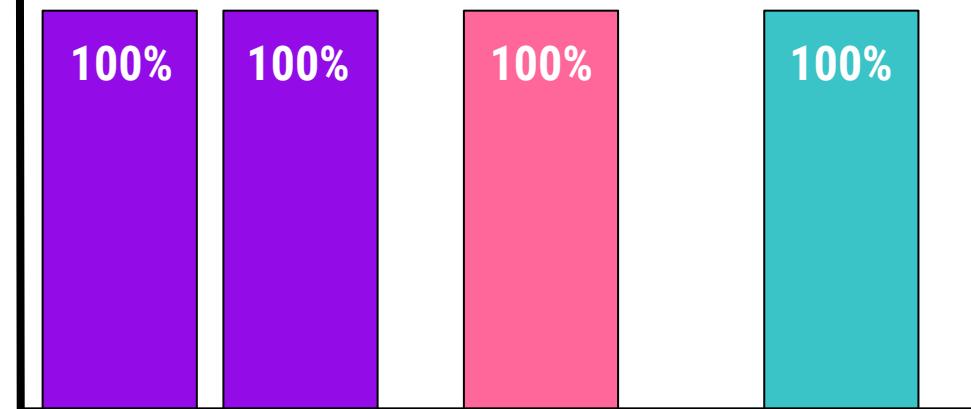
```
filter_regions()  
build_predictor()  
test_predictor()  
extract_data()  
predict_pheno()
```



phenopredict



Accuracy



functions

```
filter_regions()  
build_predictor()  
test_predictor()  
extract_data()  
predict_pheno()
```



GTEX

Genotype Tissue Expression Project
N=9,662

TCGA

The Cancer Genome Atlas
N=11,284

SRA

Sequence Read Archive
N=49,848

Training
Data

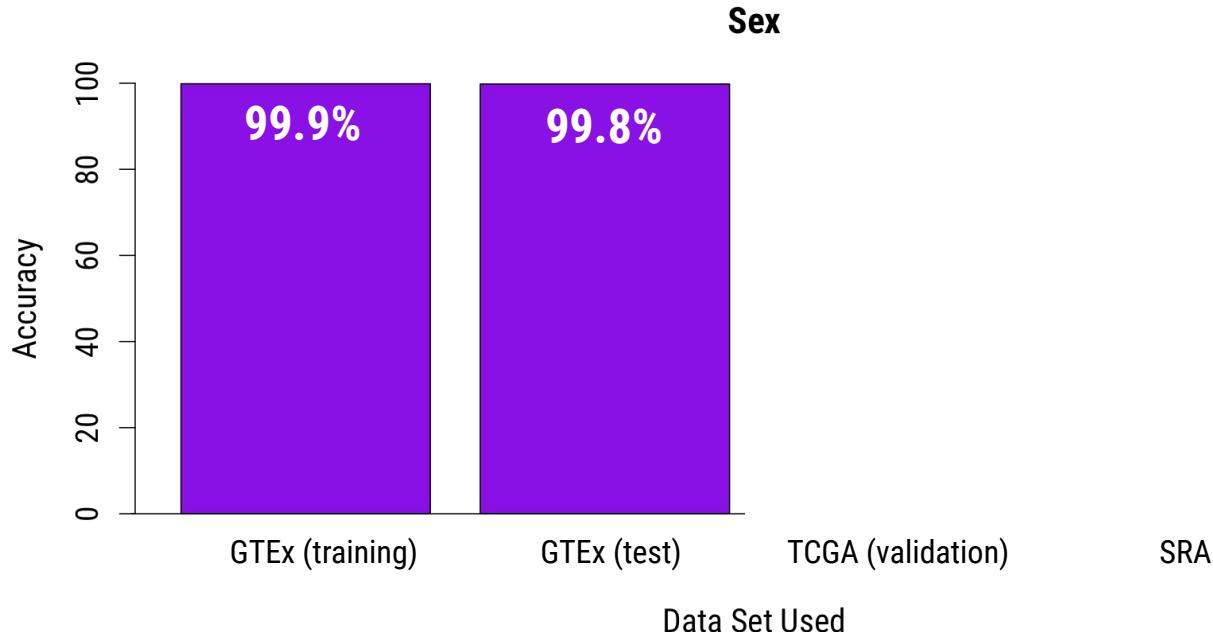
Test
Data

Validation
Data

Make
predictions!

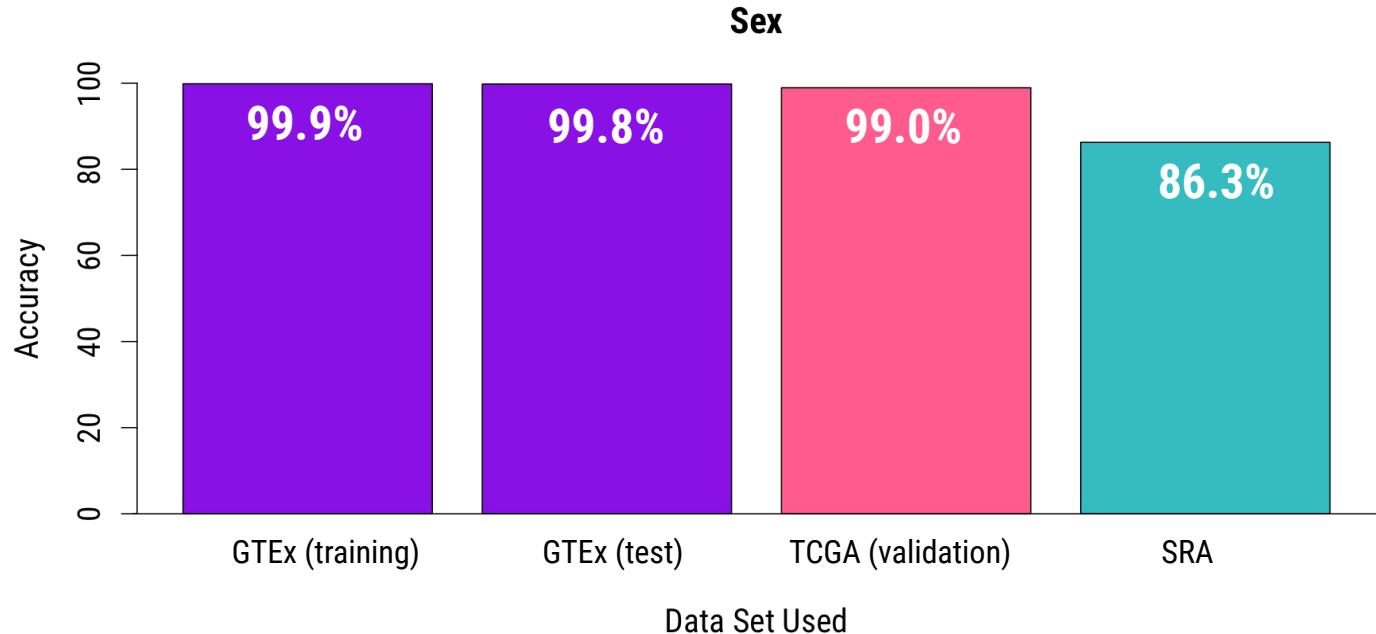
Let's get predicting...

Sex
prediction is
accurate
across data
sets



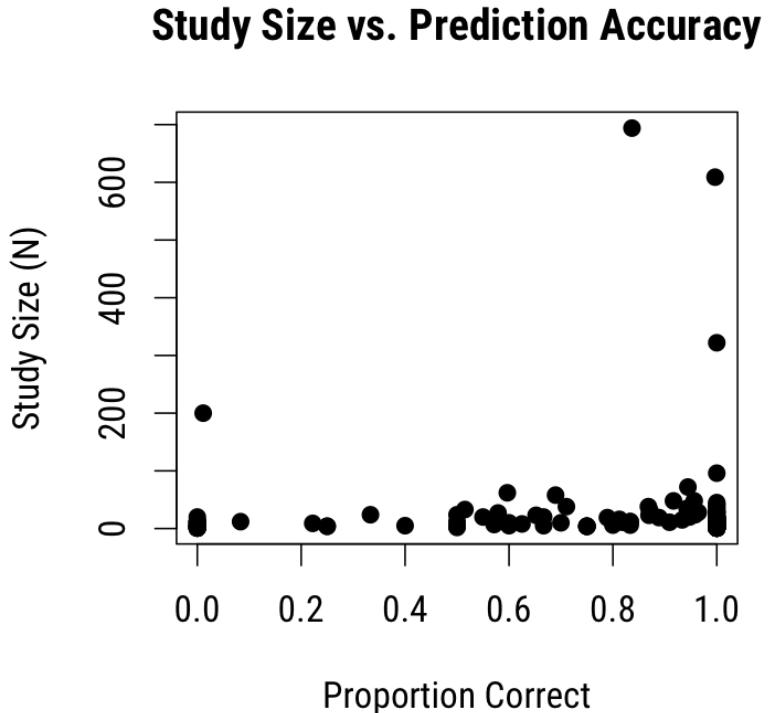
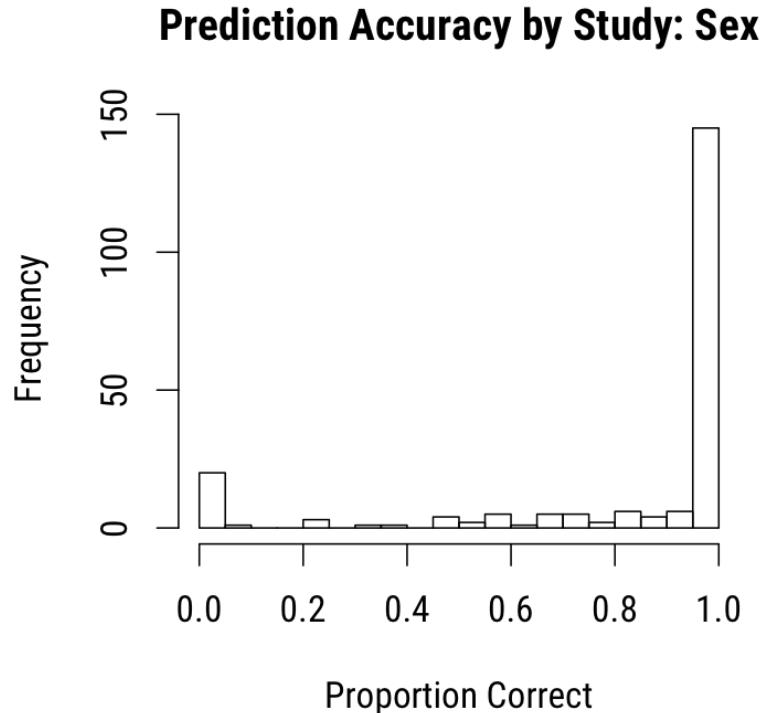
Number of Regions	40	40
Number of Samples (N)	4,769	4,769

Sex
prediction is
accurate
across data
sets

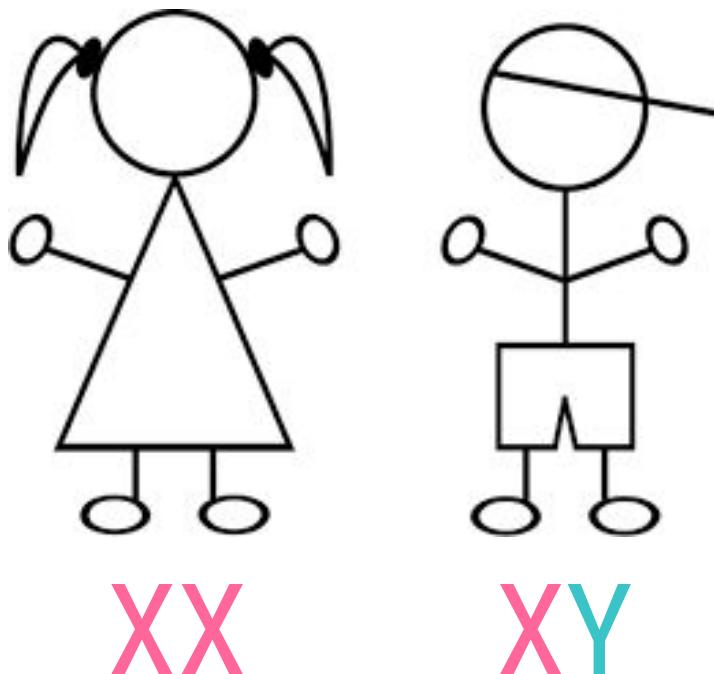


Number of Regions	40	40	40	40
Number of Samples (N)	4,769	4,769	11,245	3,640

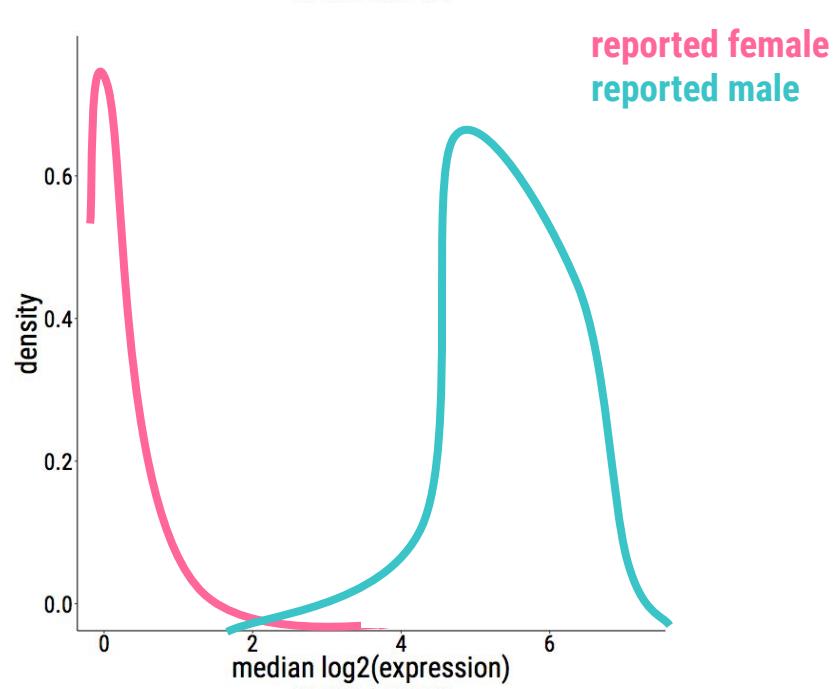
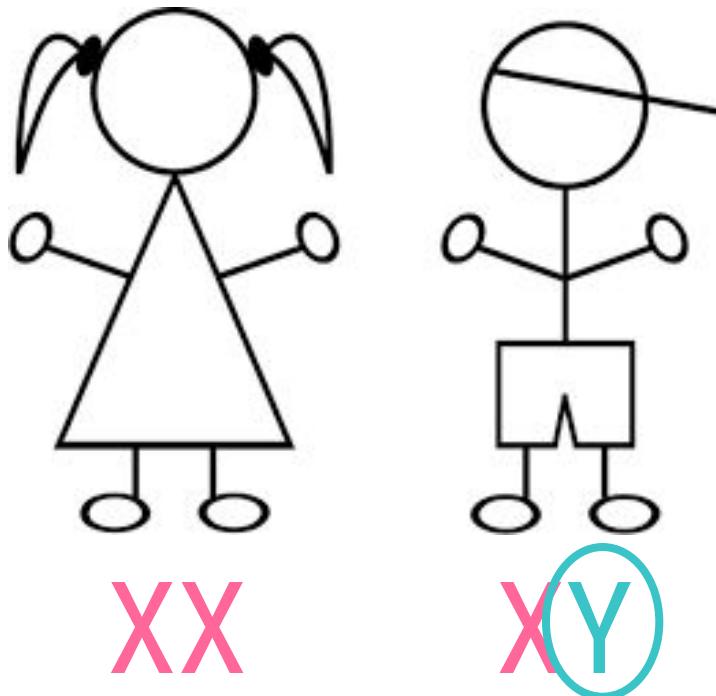
Are a few studies driving decrease in accuracy across the SRA samples?



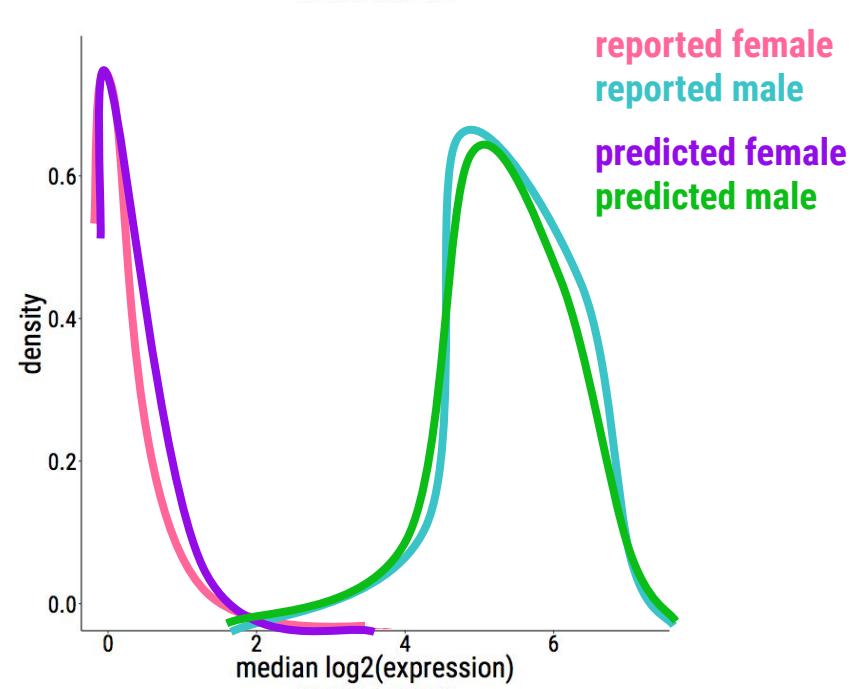
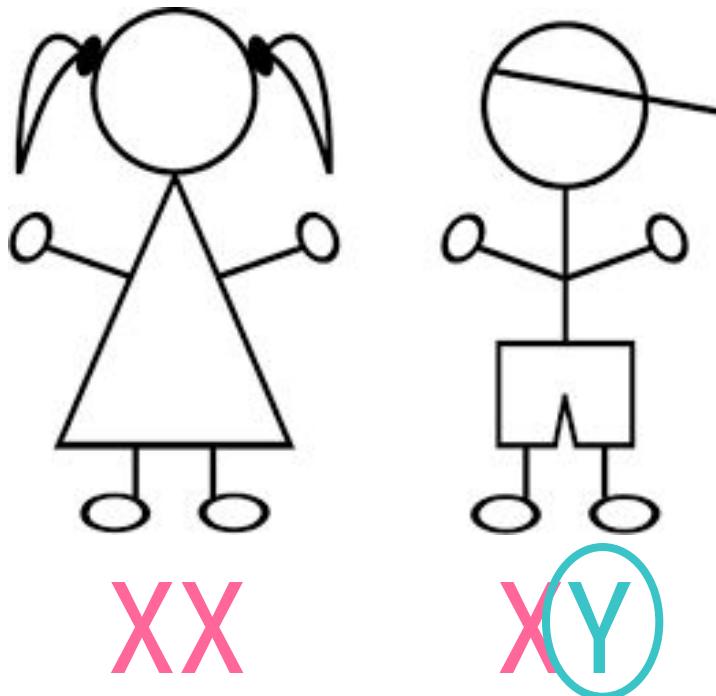
To assess misreporting of sex in the SRA, we can use Y-chromosome expression



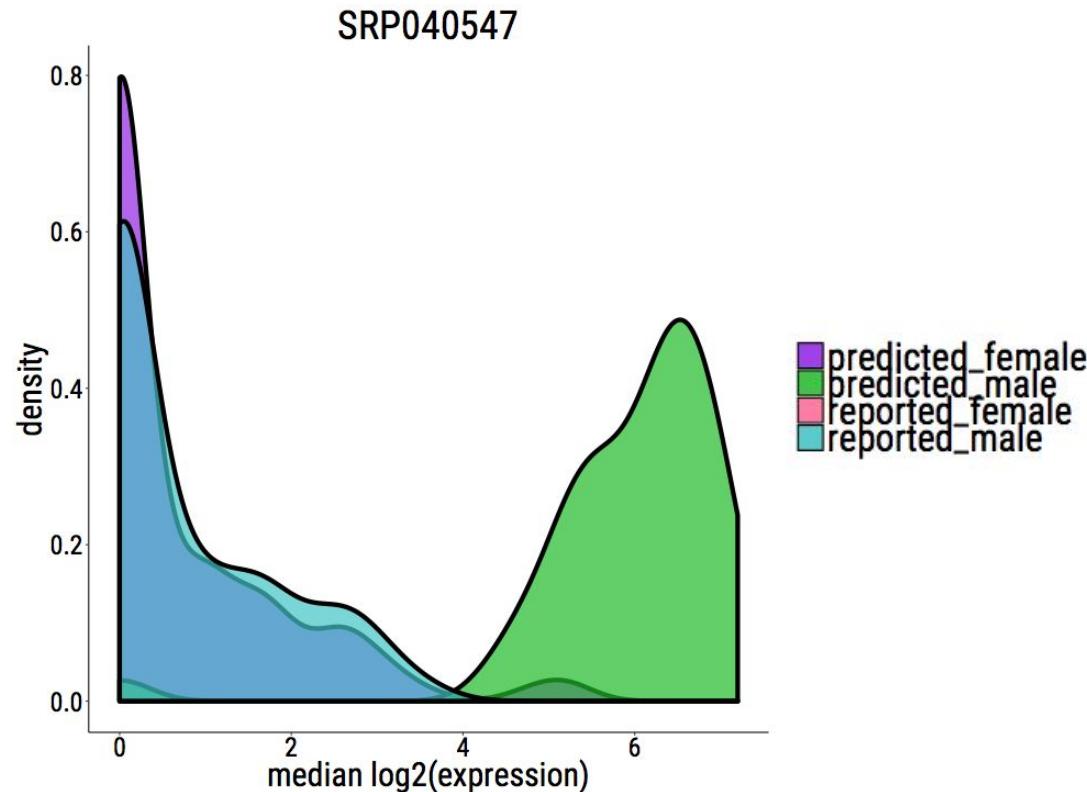
To assess misreporting of sex in the SRA, we can use Y-chromosome expression



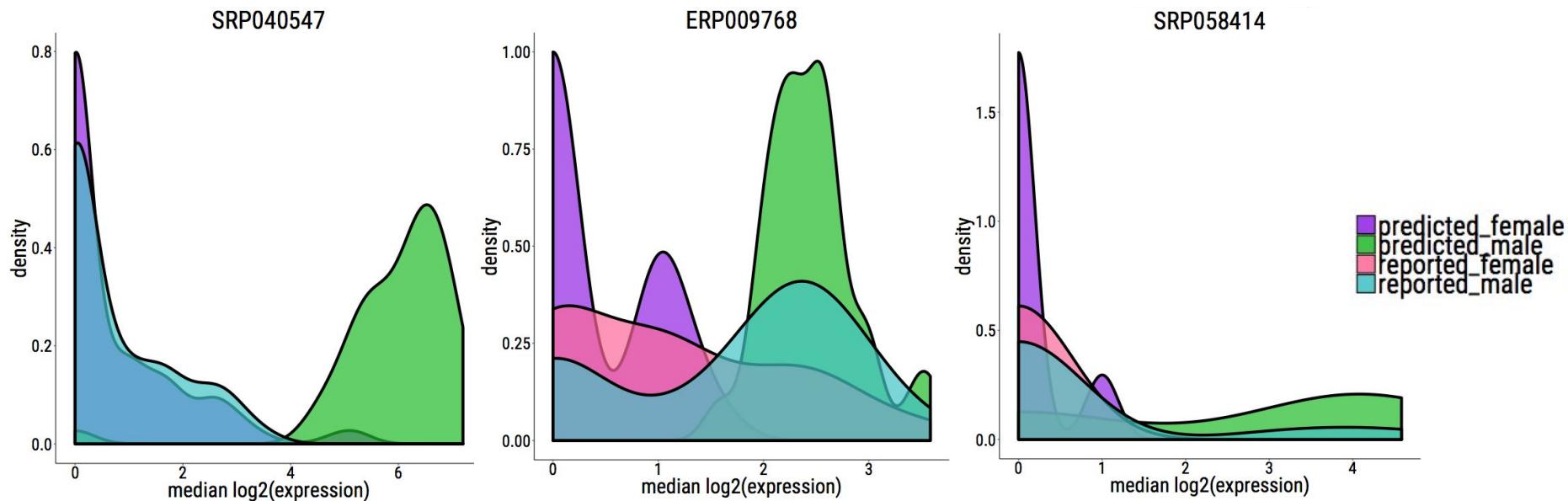
To assess misreporting of sex in the SRA, we can use Y-chromosome expression

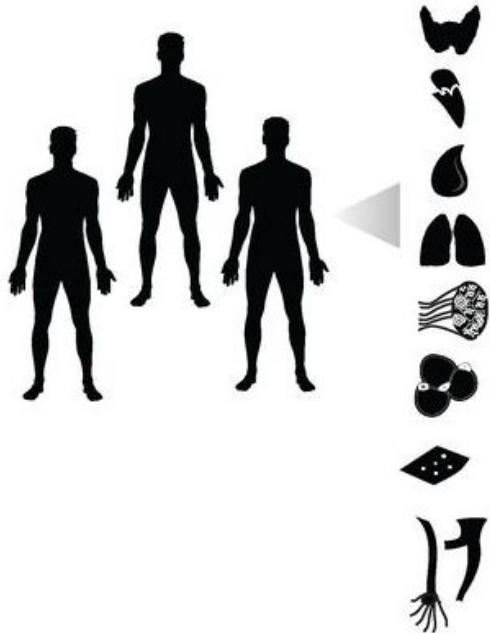


Expression from the Y chromosome suggests misreporting of sex in the SRA



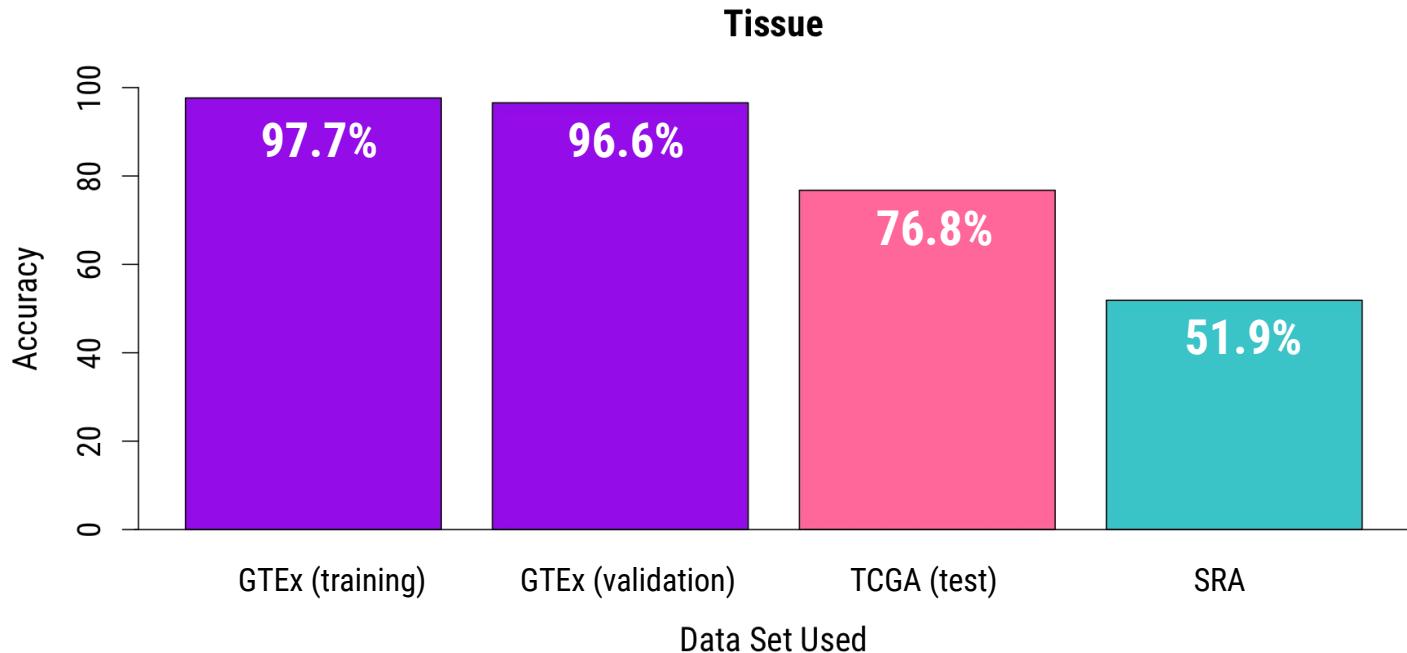
Expression from the Y chromosome suggests misreporting of sex in the SRA





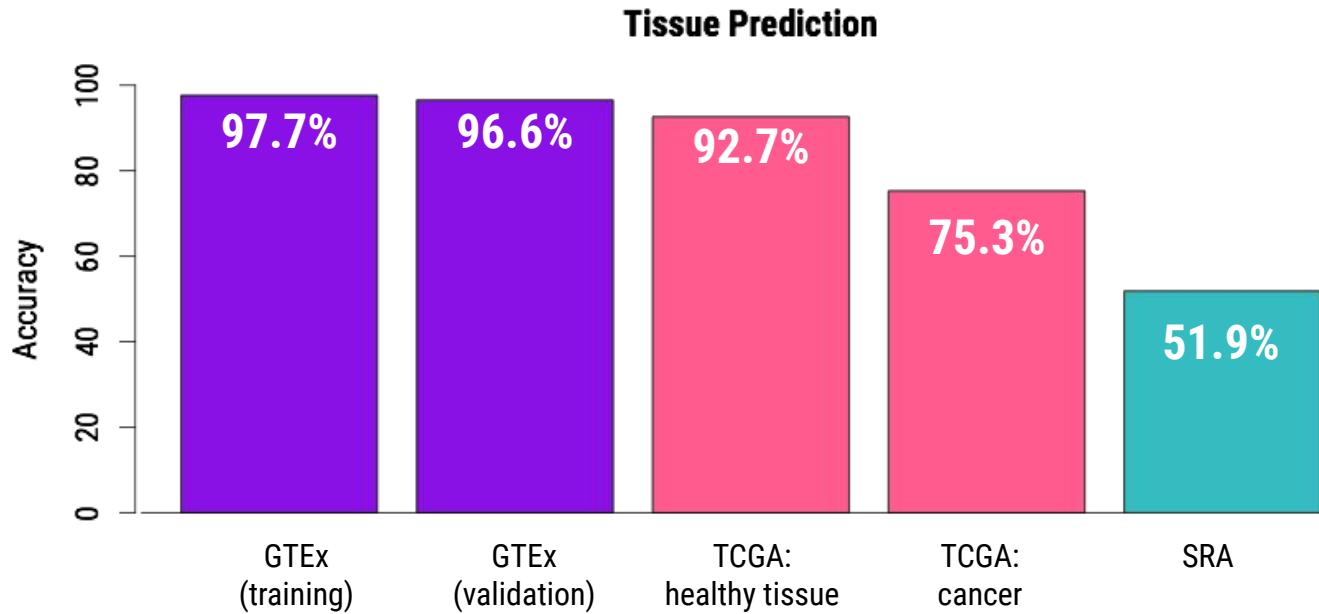
Can we use
expression data
to predict
tissue?

Tissue
prediction is
accurate
across data
sets



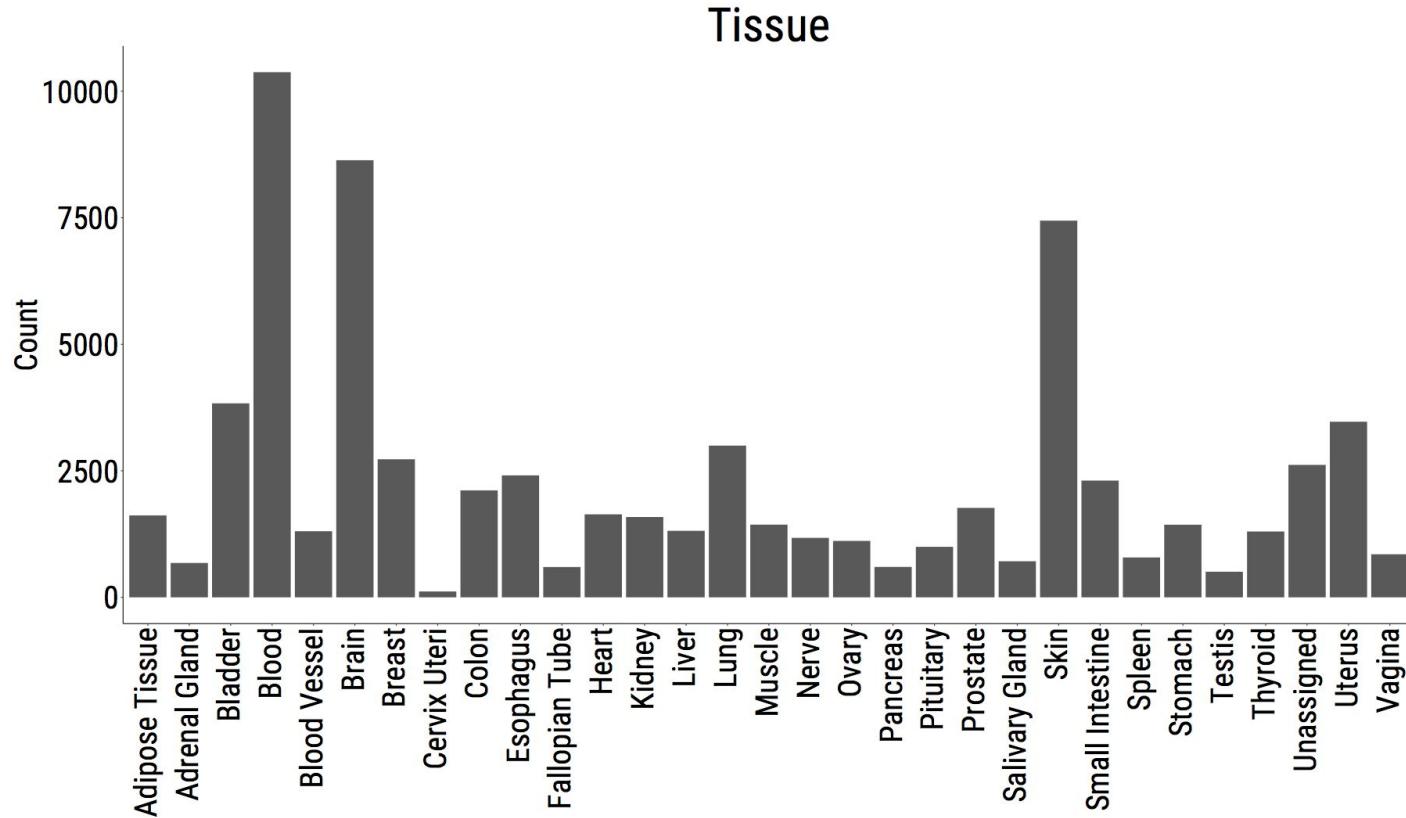
Number of Regions	2,281	2,281	2,281	2,281
Number of Samples (N)	4,769	4,769	7,317	8,951

Prediction is
more
accurate in
healthy
tissue



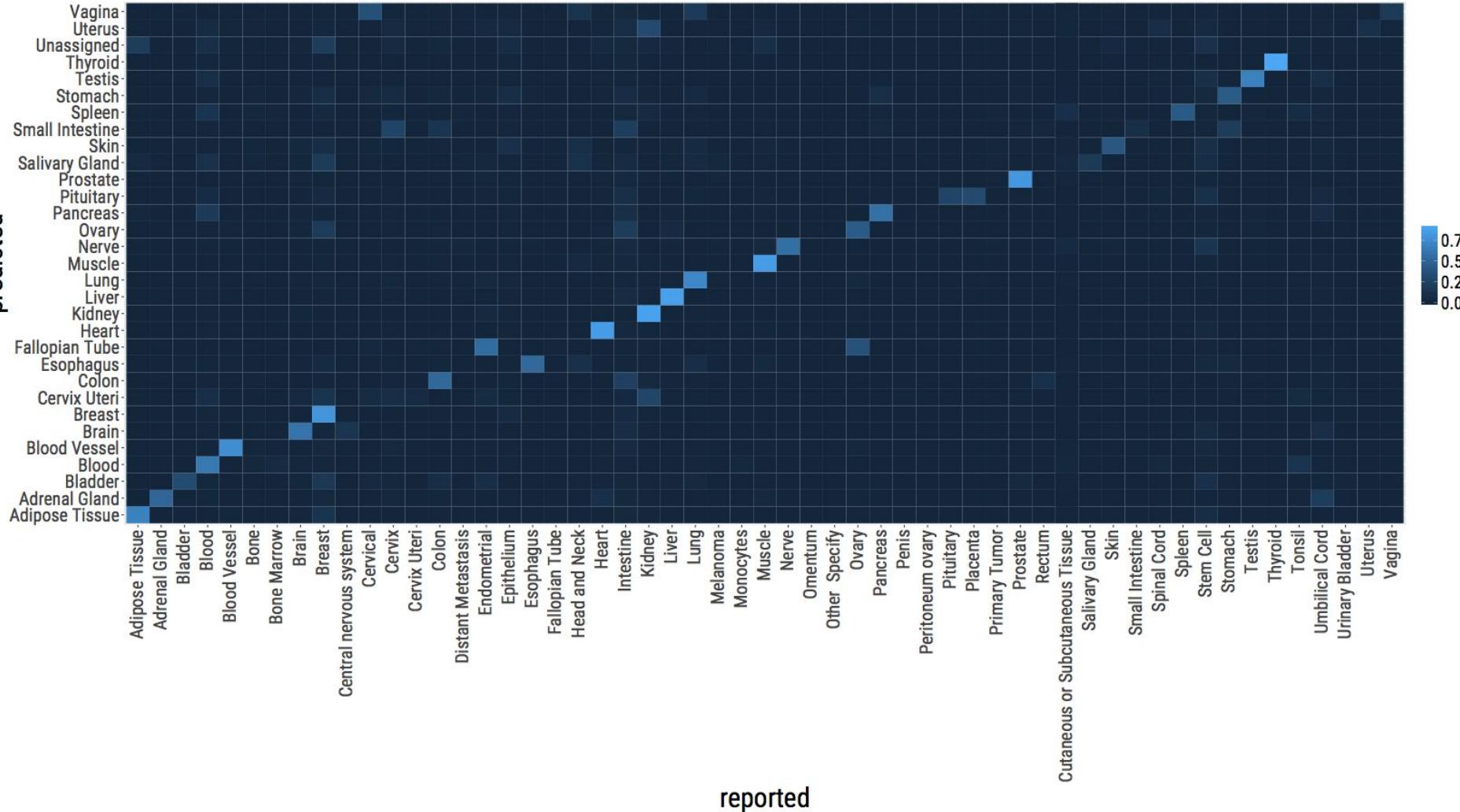
Number of Regions	2,281	2,281	2,281	2,281	2,281
Number of Samples (N)	4,769	4,769	613	6,704	8,951

Across the samples in *recount*, brain, blood, and skin are the three most frequently predicted tissues types



Tissue

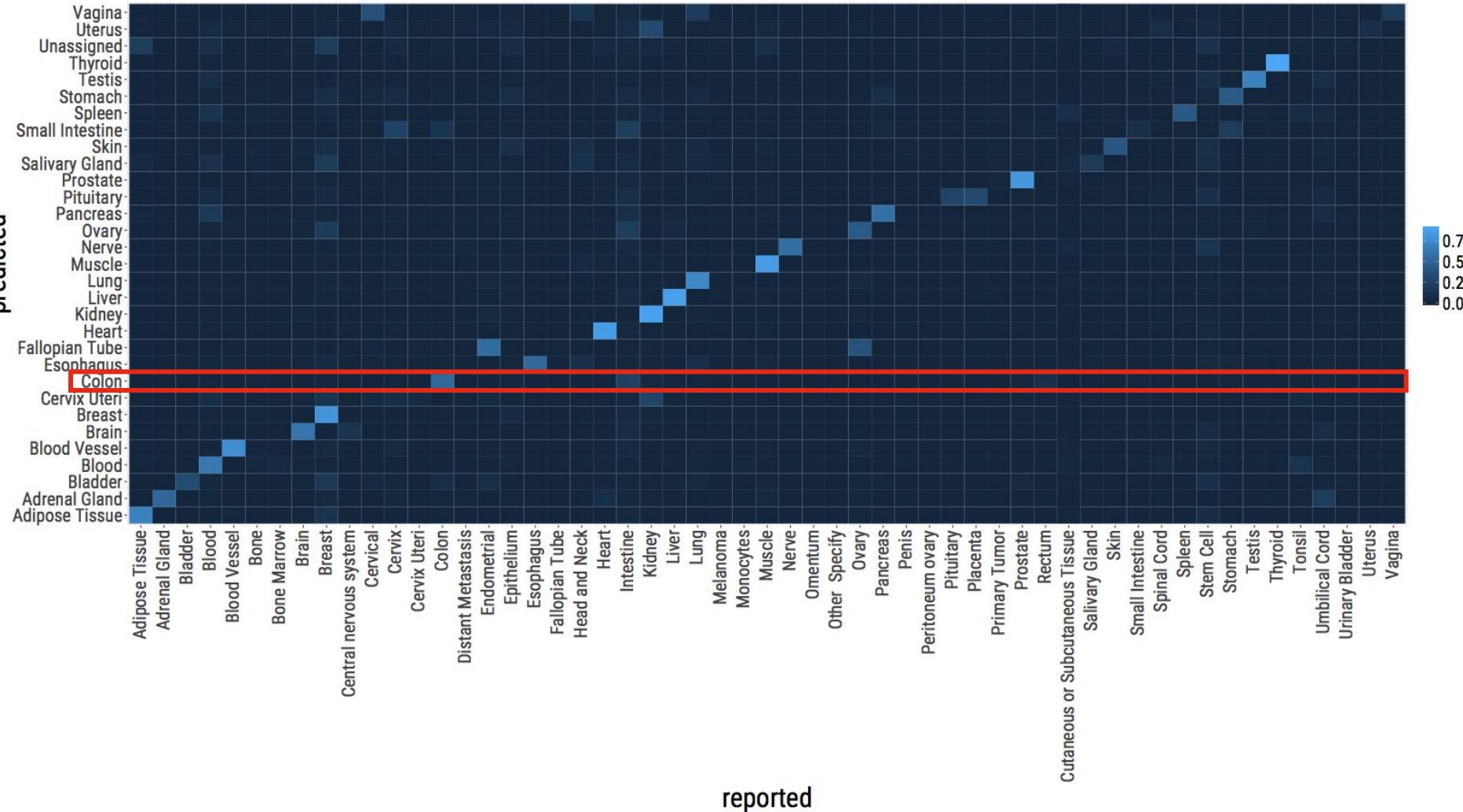
predicted



reported

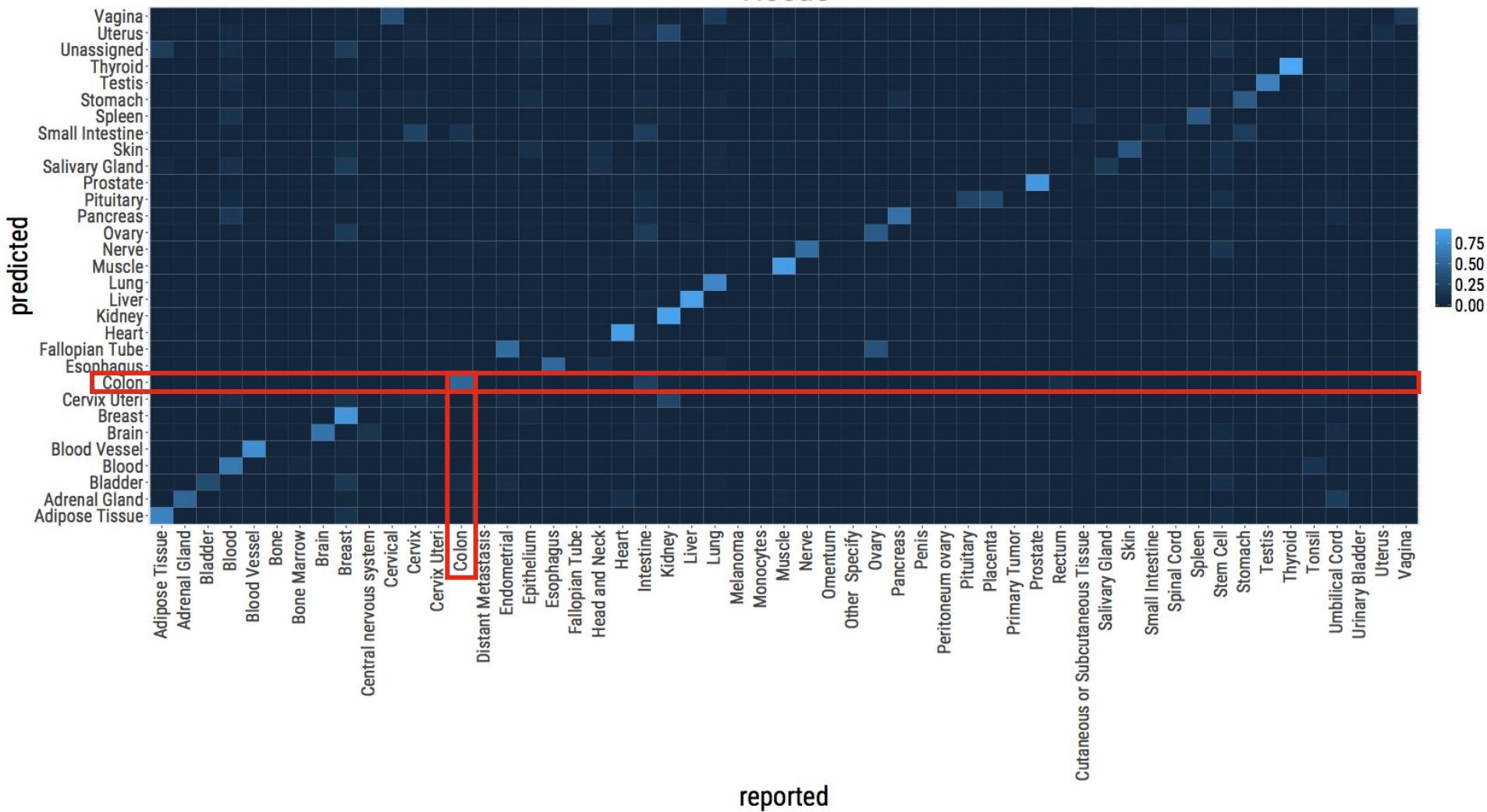
Tissue

predicted



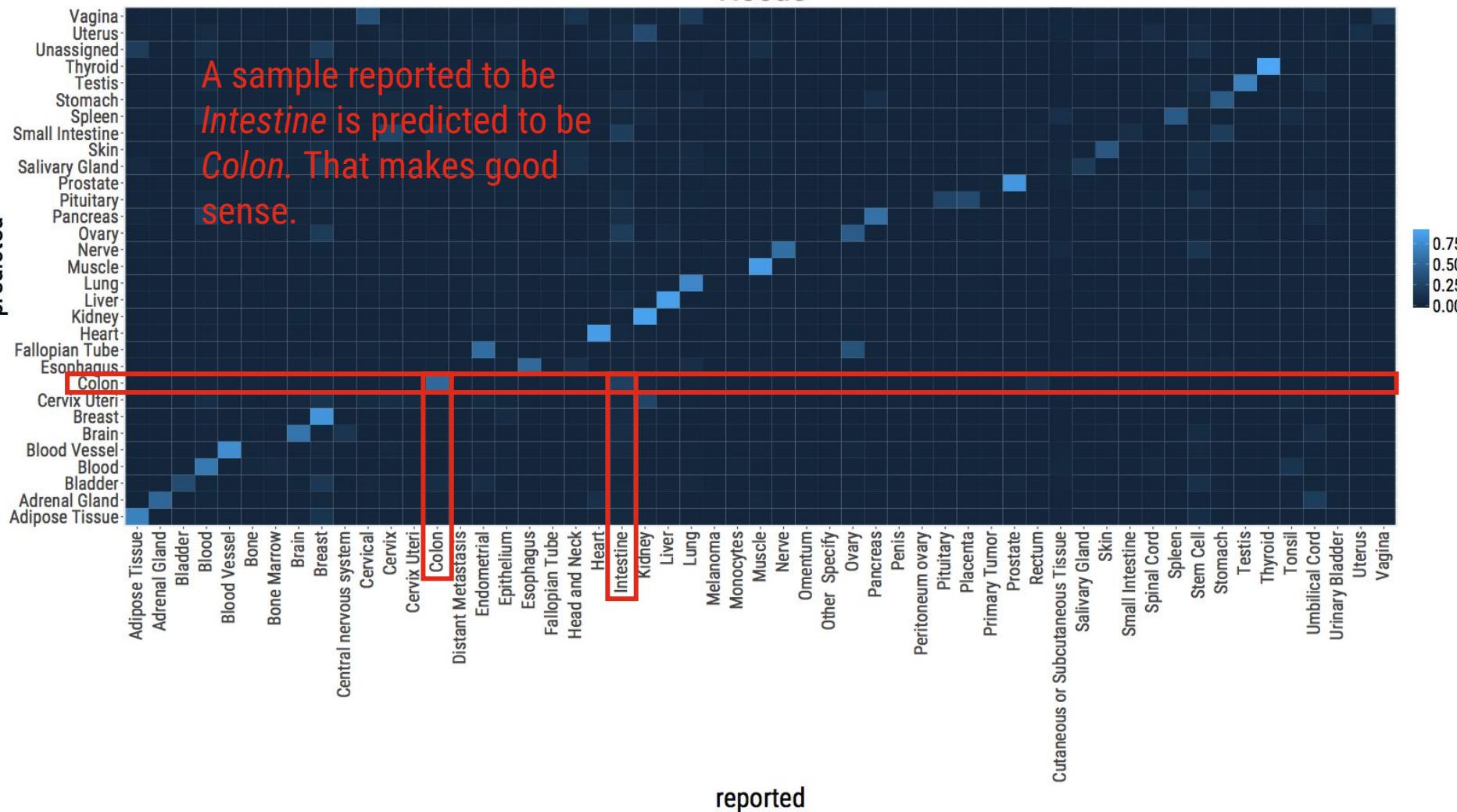
reported

Tissue

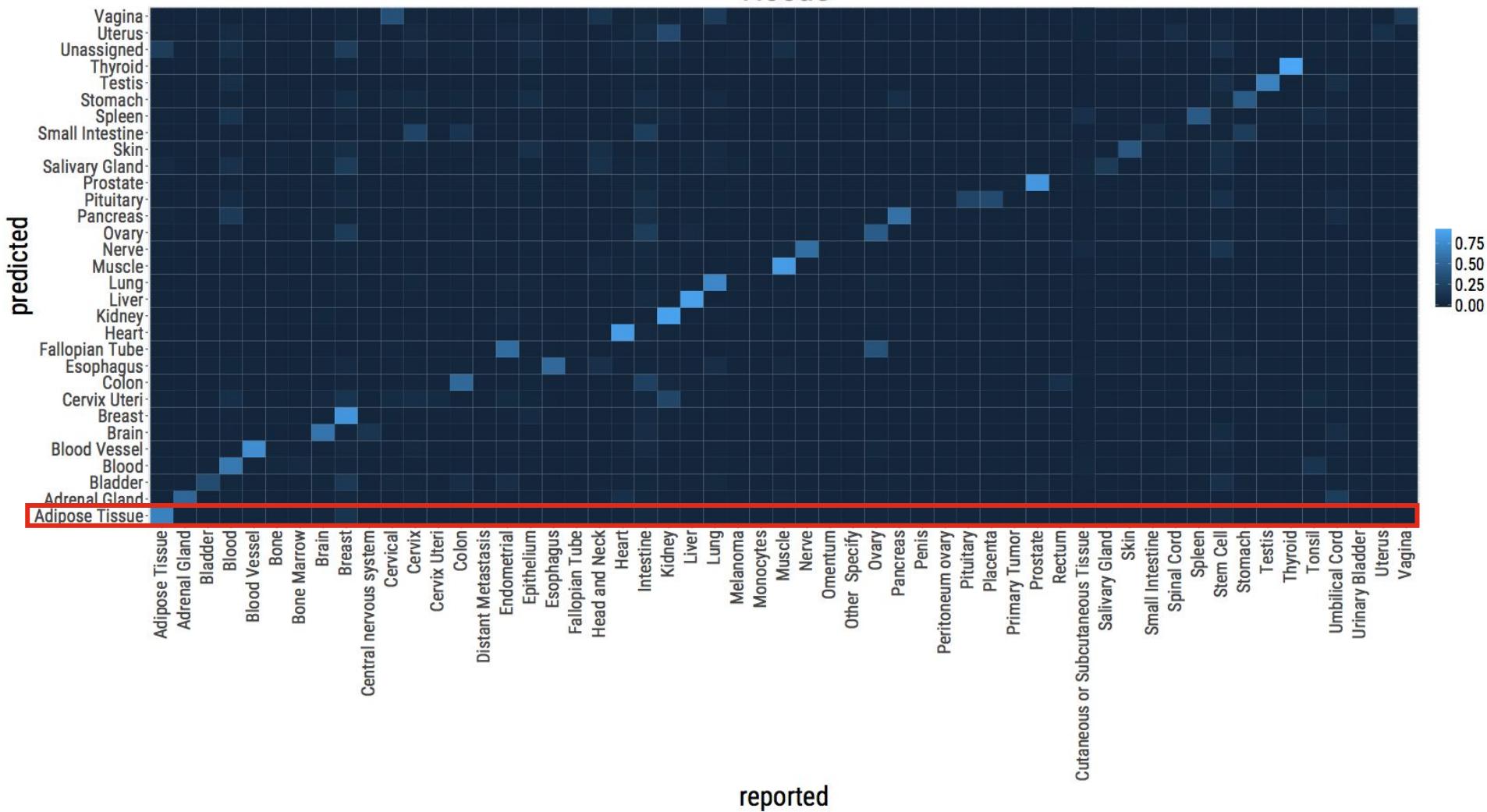


Tissue

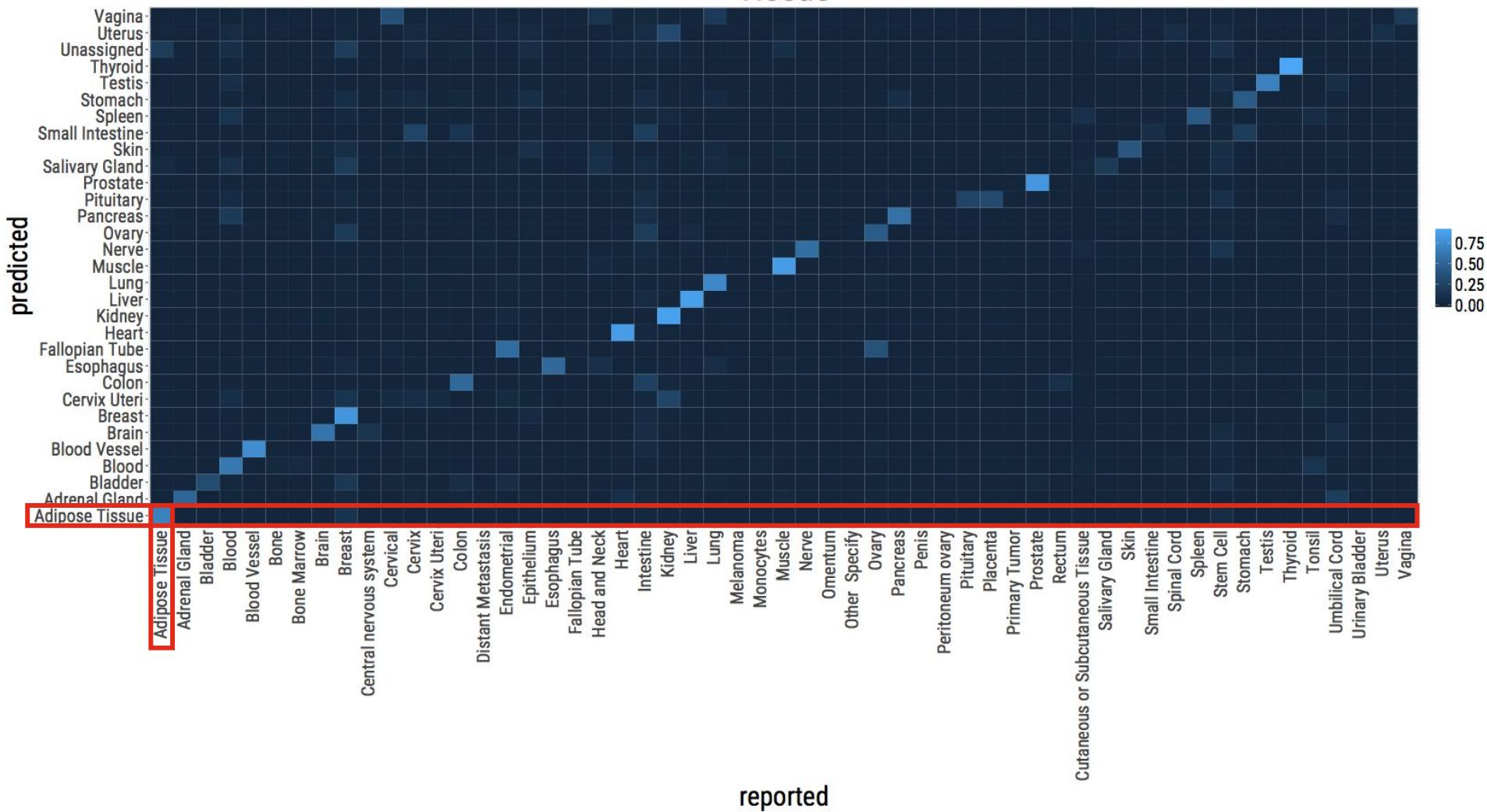
A sample reported to be *Intestine* is predicted to be *Colon*. That makes good sense.



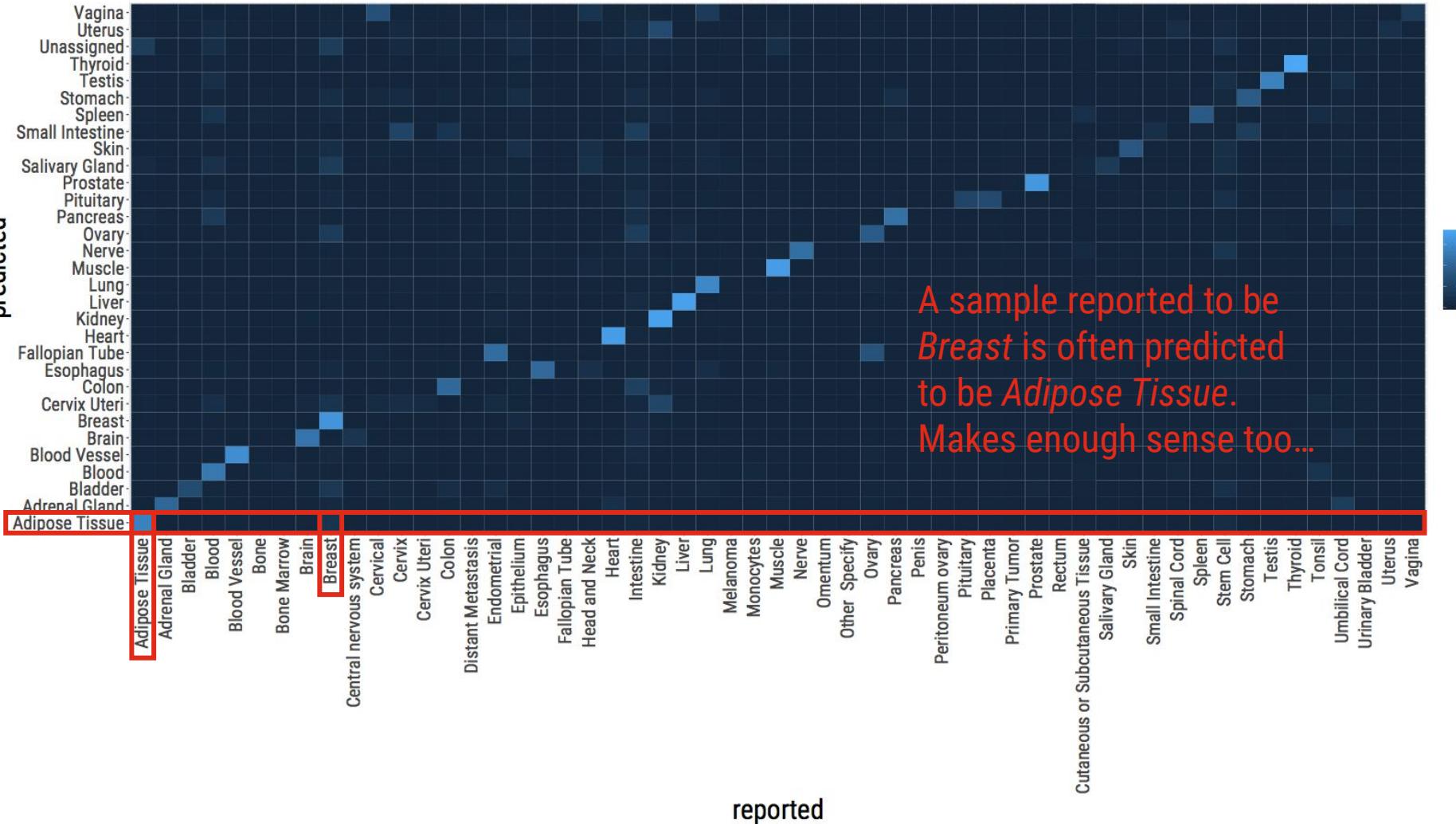
Tissue



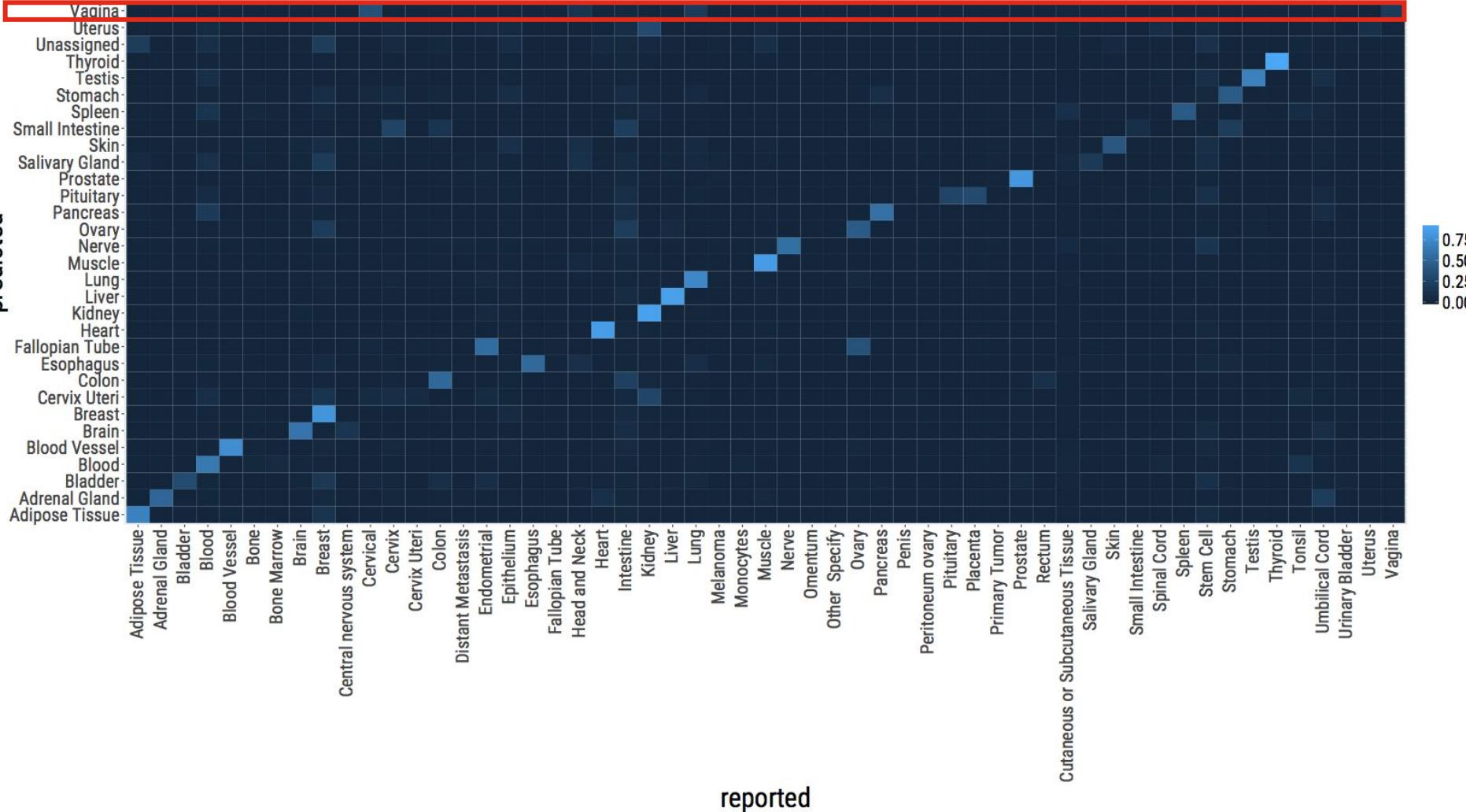
Tissue



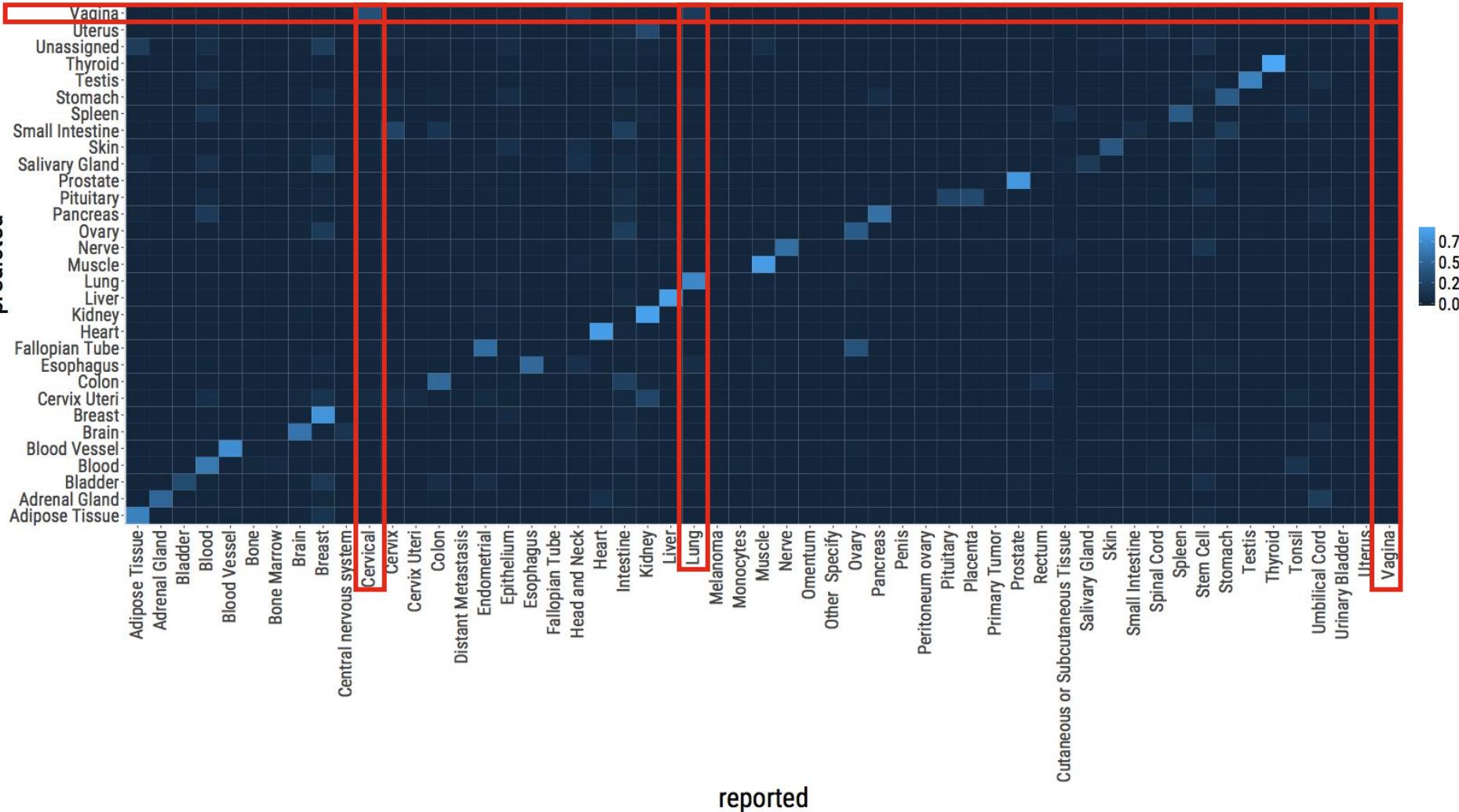
Tissue



Tissue



Tissue



Tissue

predicted

Vagina
Uterus
Unassigned
Thyroid
Testis
Stomach
Spleen
Small Intestine
Skin
Salivary Gland
Prostate
Pituitary
Pancreas
Ovary
Nerve
Muscle
Lung
Liver
Kidney
Heart
Fallopian Tube
Esophagus
Colon
Cervix Uteri
Breast
Brain
Blood Vessel
Blood
Bladder
Adrenal Gland
Adipose Tissue

Adipose Tissue

Adrenal Gland

Bladder

Blood

Blood Vessel

Bone

Bone Marrow

Brain

Breast

Central nervous system

Cervical

Cervix

Cervix Uteri

Colon

Distant Metastasis

Endometrial

Epithelium

Esophagus

Fallopian Tube

Head and Neck

Heart

Intestine

Kidney

Liver

Lung

Melanoma

Monocytes

Muscle

Nerve

Omentum

Other Specify

Ovary

Pancreas

Penis

Peritoneum ovary

Pituitary

Placenta

Primary Tumor

Prostate

Rectum

Cutaneous or Subcutaneous Tissue

Salivary Gland

Skin

Small Intestine

Spinal Cord

Spleen

Stem Cell

Stomach

Testis

Thyroid

Tonsil

Umbilical Cord

Urinary Bladder

Vagina

reported

-＼(ツ)／-

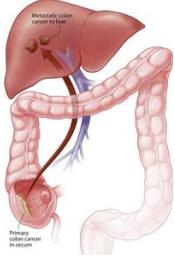
But sometimes the predictions and reported tissues make less sense...

Tissue prediction is largely accurate across *recount2*

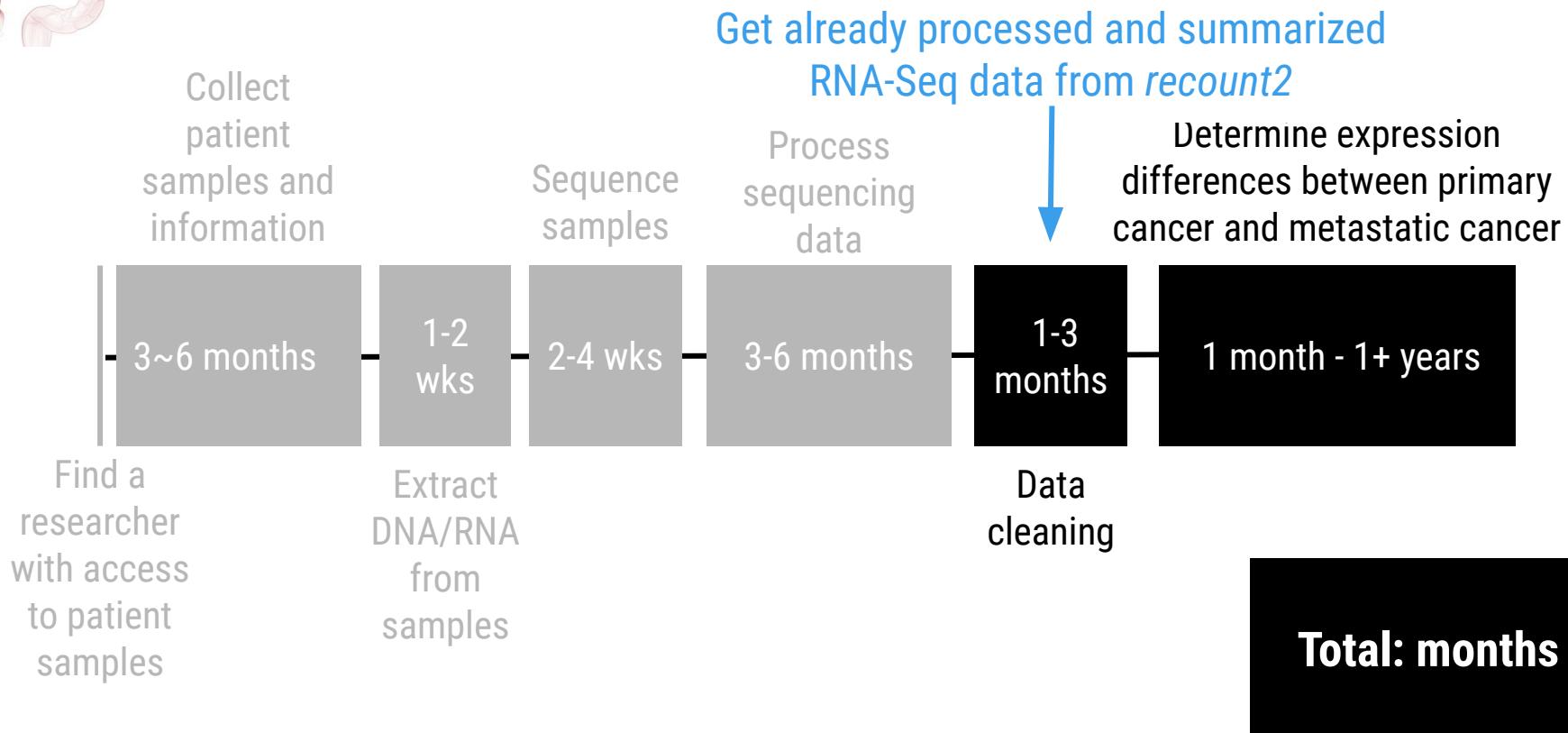
Tissue can be accurately predicted from expression data.

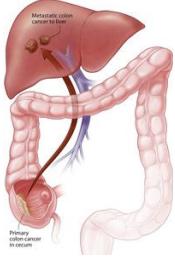
Discordant predictions are often made to biologically similar tissues.

Sometimes, predictions are inaccurate.

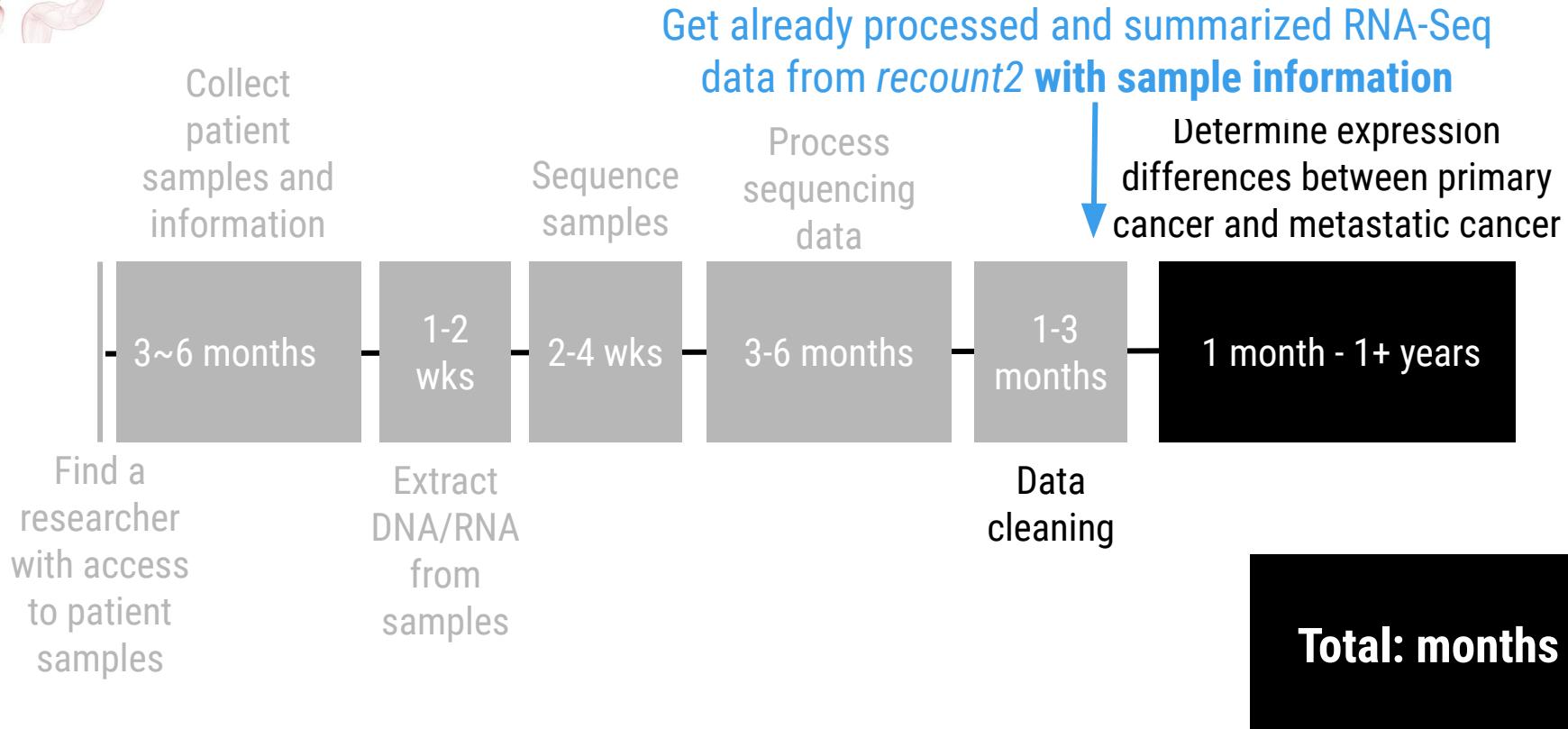


What makes primary cancer different than metastatic cancer?





What makes primary cancer different than metastatic cancer?

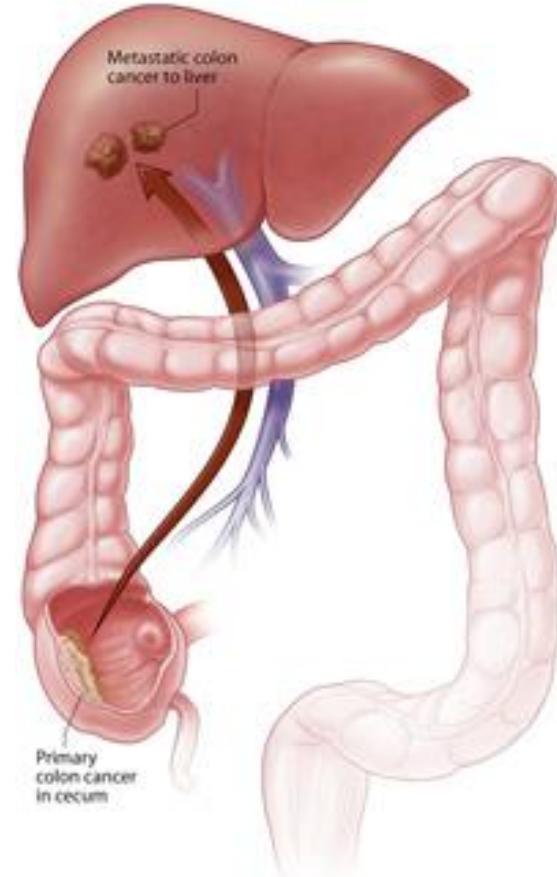


Project Goal: Take publicly available RNA-Seq data and make it available and easy-to-use

- ☒ Comprehensive
- ☒ Easy to Get
- ☒ Useful for future study

Ok. Ok. What about
actually *using* all of
these predictions...?

What makes primary cancer different than metastatic cancer?



A nineteen gene-based risk score classifier predicts prognosis of colorectal cancer patients



CrossMark

Seon-Kyu Kim^{a,1}, Seon-Young Kim^{a,1}, Jeong-Hwan Kim^a, Seon Ae Roh^{b,c},
Dong-Hyung Cho^{c,d}, Yong Sung Kim^{a,c,**}, Jin Cheon Kim^{b,c,*}

^aMedical Genomics Research Centre, Korea Research Institute of Bioscience and Biotechnology, Daejeon, Korea

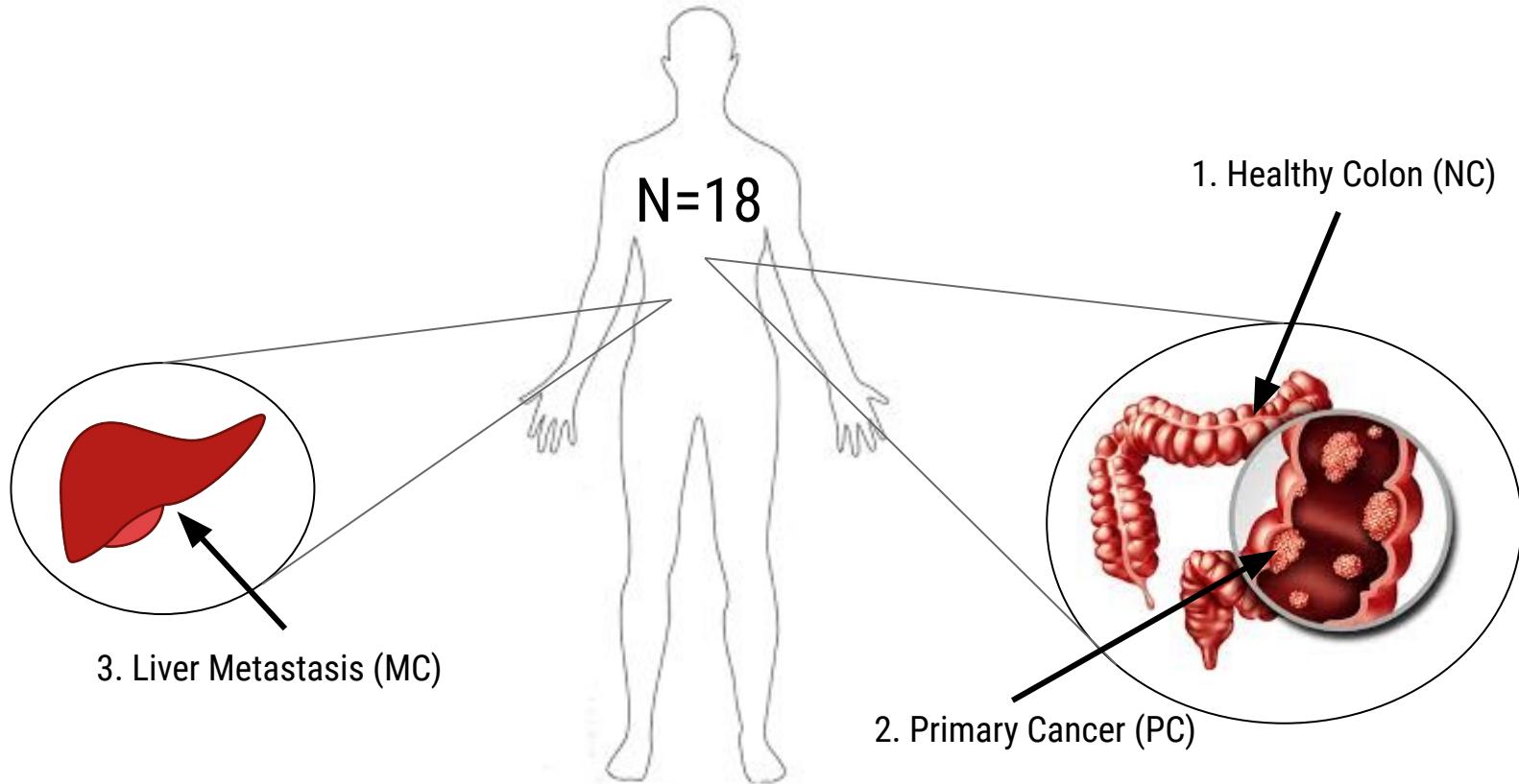
^bDepartment of Surgery, University of Ulsan College of Medicine, Seoul, Korea

^cDepartment of Cancer Research, Institute of Innovative Cancer Research and Asan Institute for Life Sciences,
Asan Medical Centre, Seoul, Korea

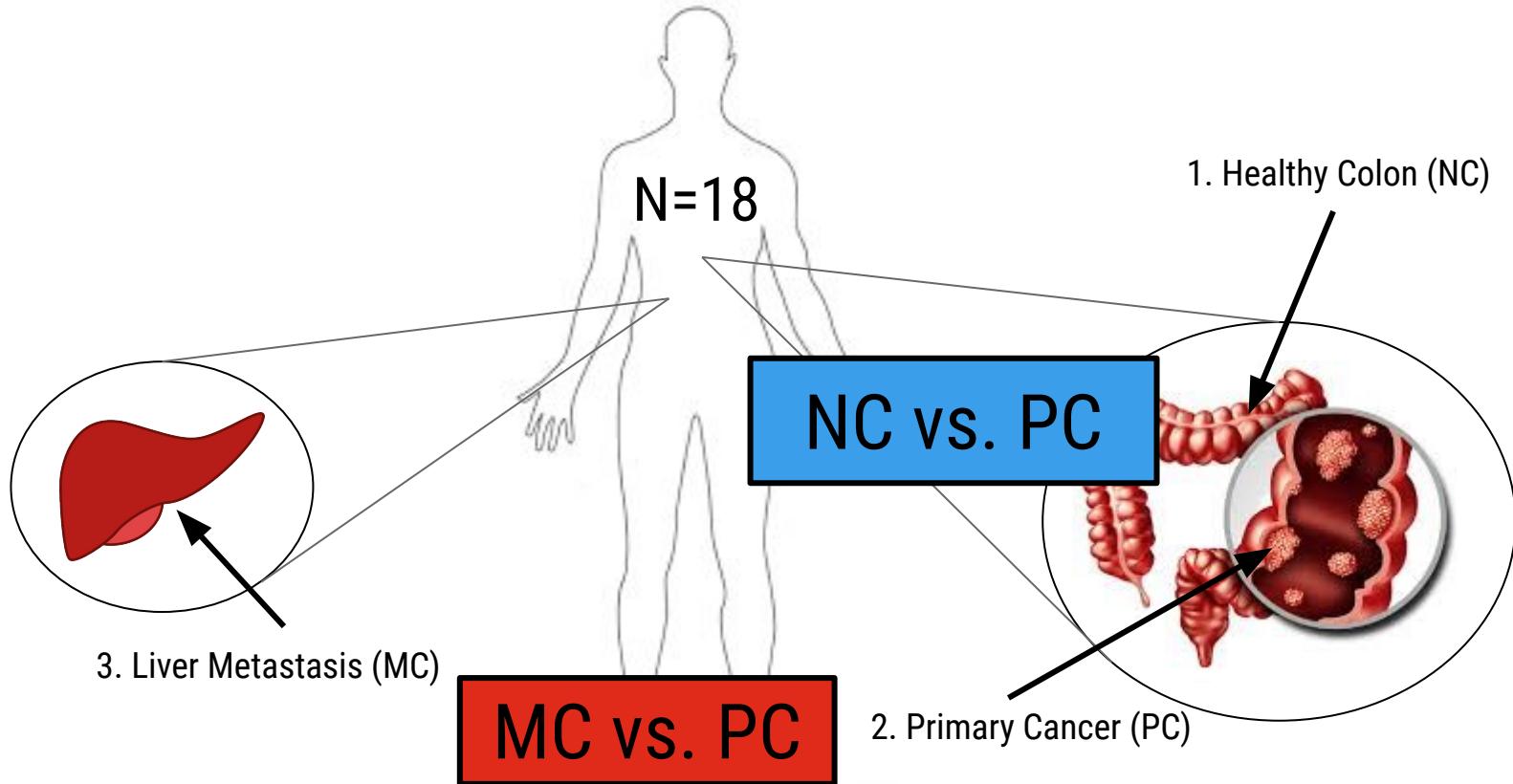
^dGraduate School of East-West Medical Science, Kyung Hee University, Gyeonggi-do, Korea

Molecular Oncology, July 2014

Kim et al. analysis looked to identify genes that contribute to metastasis in colon cancer.



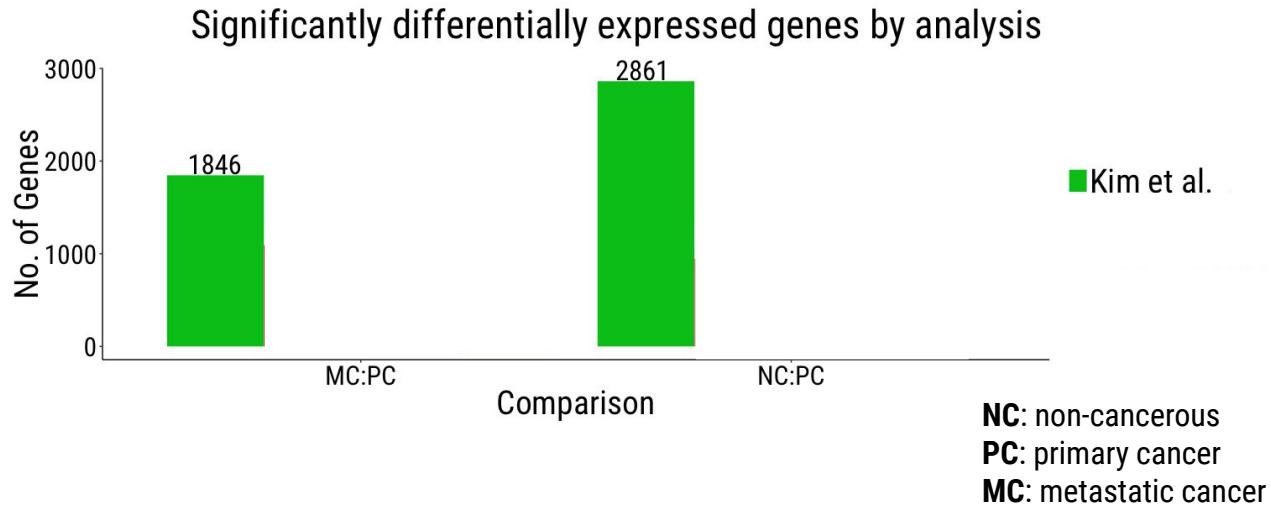
Kim et al. analysis looked to identify genes that contribute to metastasis in colon cancer.



Predictions can
be used to:

(1) Identify
studies of
interest

(2) appropriately
analyze data

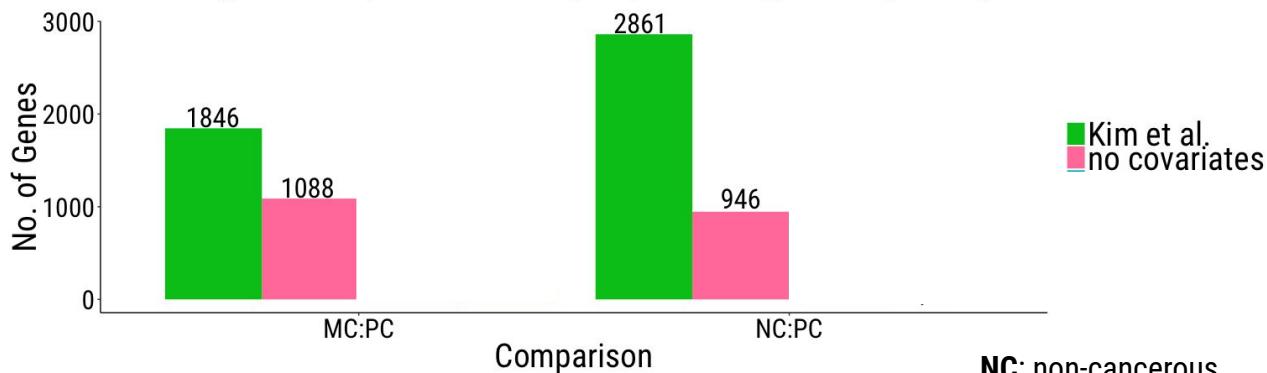


Predictions can
be used to:

(1) Identify
studies of
interest

(2) appropriately
analyze data

Significantly differentially expressed genes by analysis

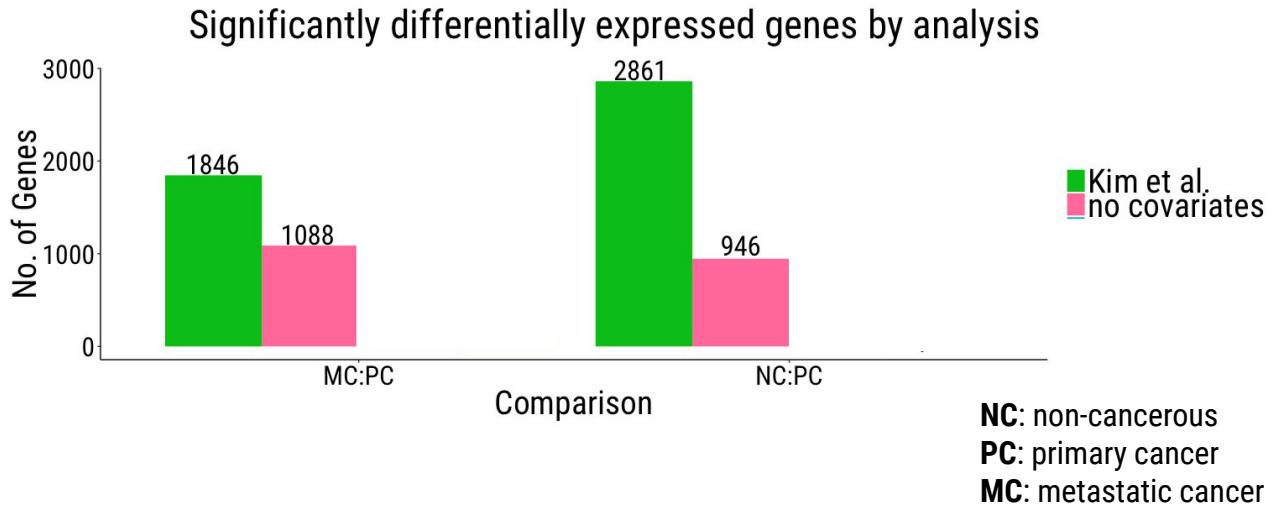


NC: non-cancerous
PC: primary cancer
MC: metastatic cancer

Predictions can
be used to:

(1) Identify
studies of
interest

(2) appropriately
analyze data

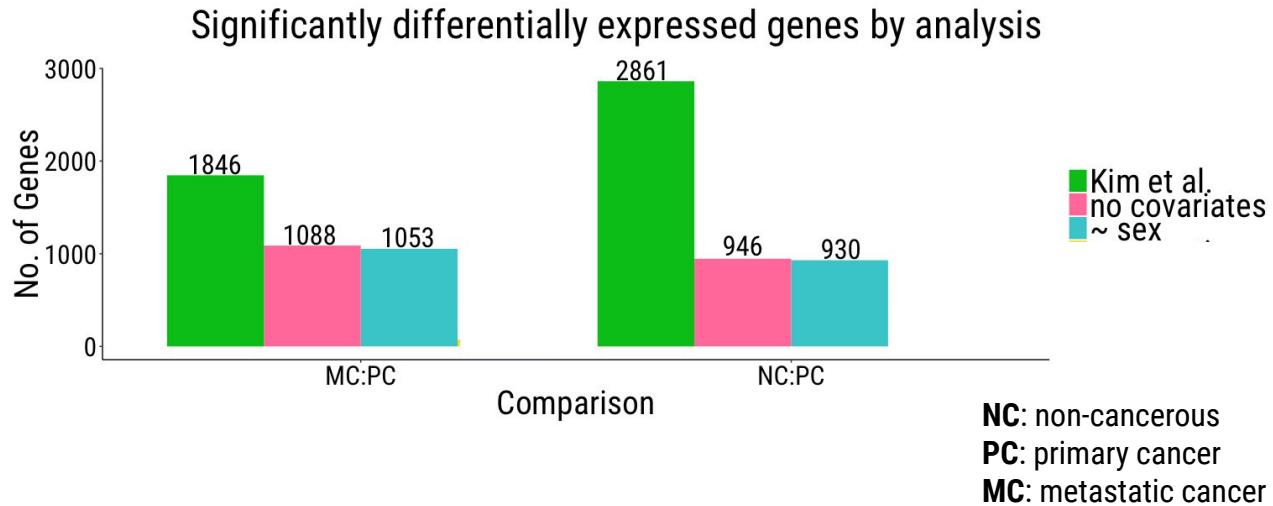


Are the same genes
found when sex is
included in the analysis?

Predictions can
be used to:

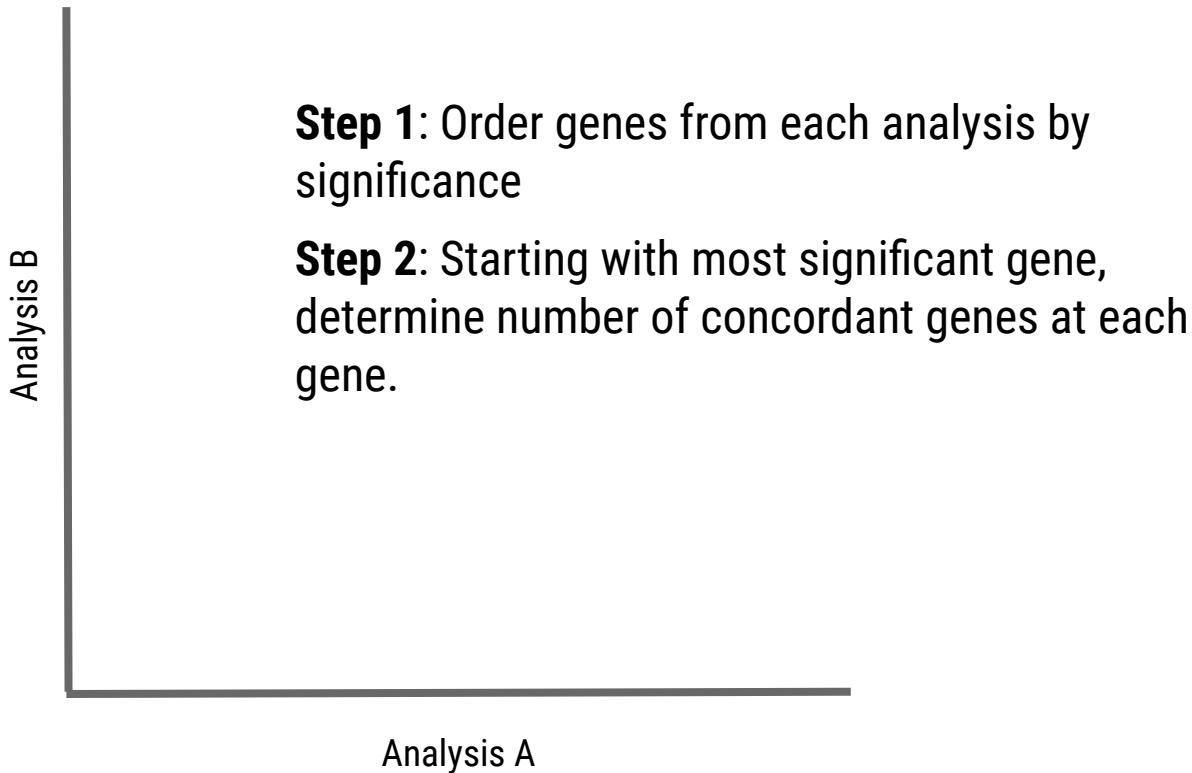
(1) Identify
studies of
interest

(2) appropriately
analyze data



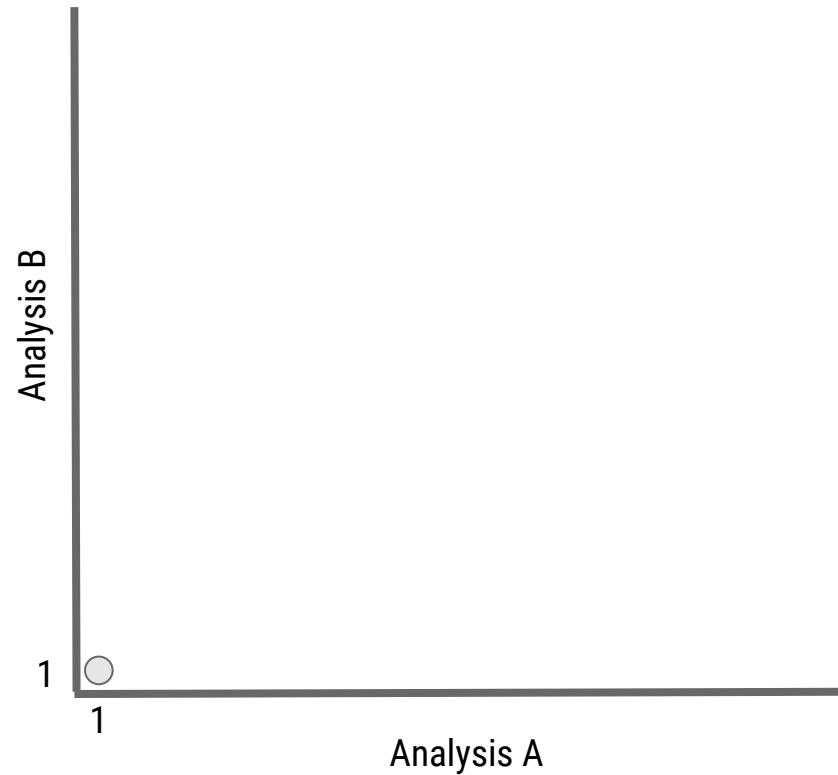
Concordance at the Top (CAT) Plots

How similar are the results from Analysis A and Analysis B ?



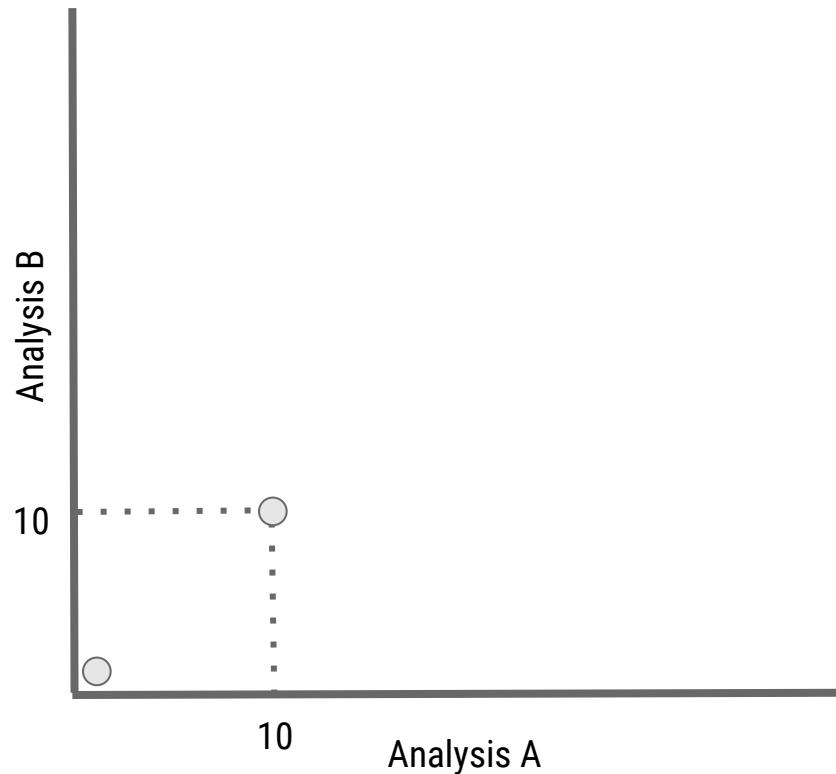
Concordance at the Top (CAT) Plots

How similar are the results from Analysis A and Analysis B ?



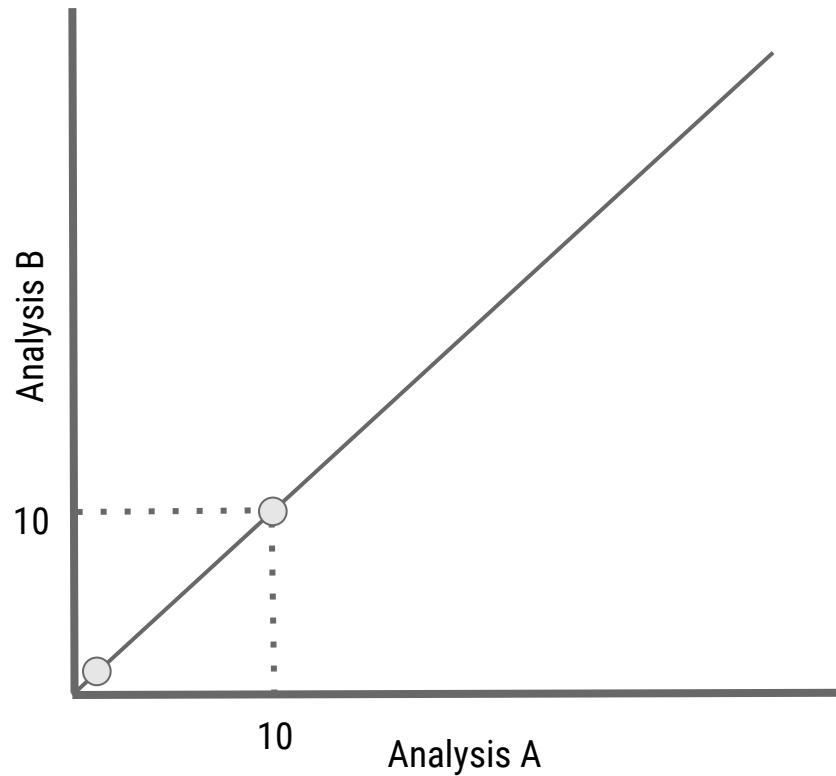
Concordance at the Top (CAT) Plots

How similar are the results from Analysis A and Analysis B ?



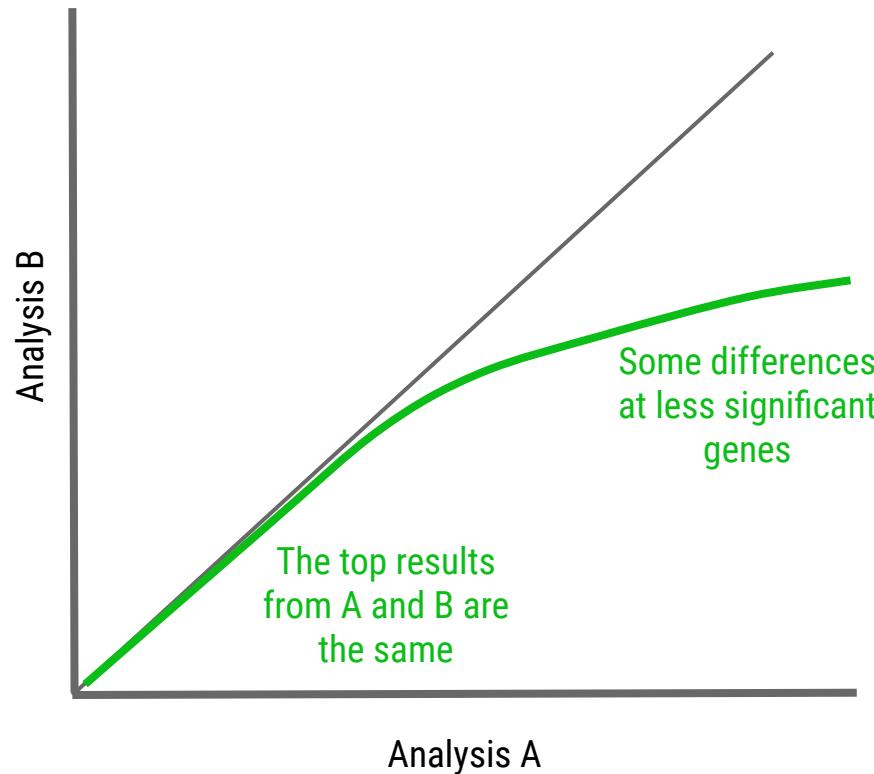
Concordance at the Top (CAT) Plots

How similar are the results from Analysis A and Analysis B ?



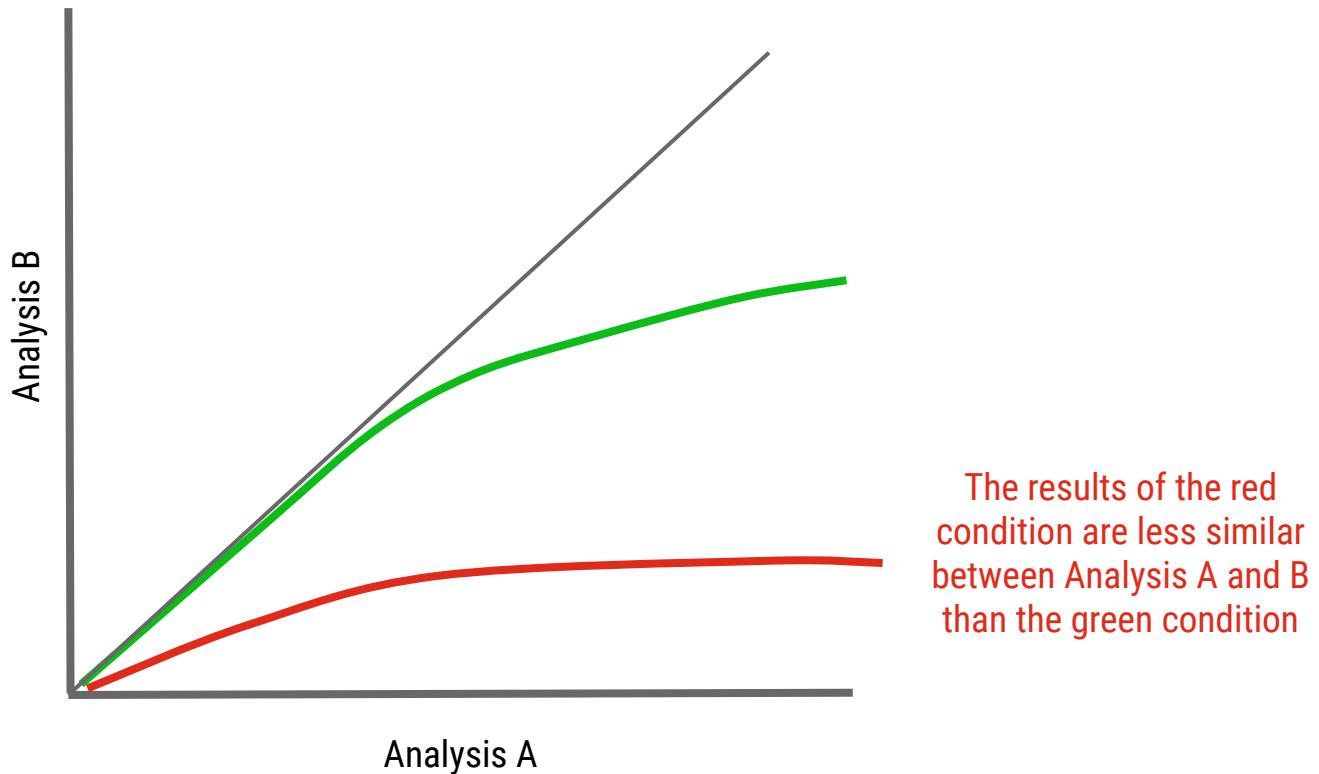
Concordance at the Top (CAT) Plots

How similar are the results from Analysis A and Analysis B ?



Concordance at the Top (CAT) Plots

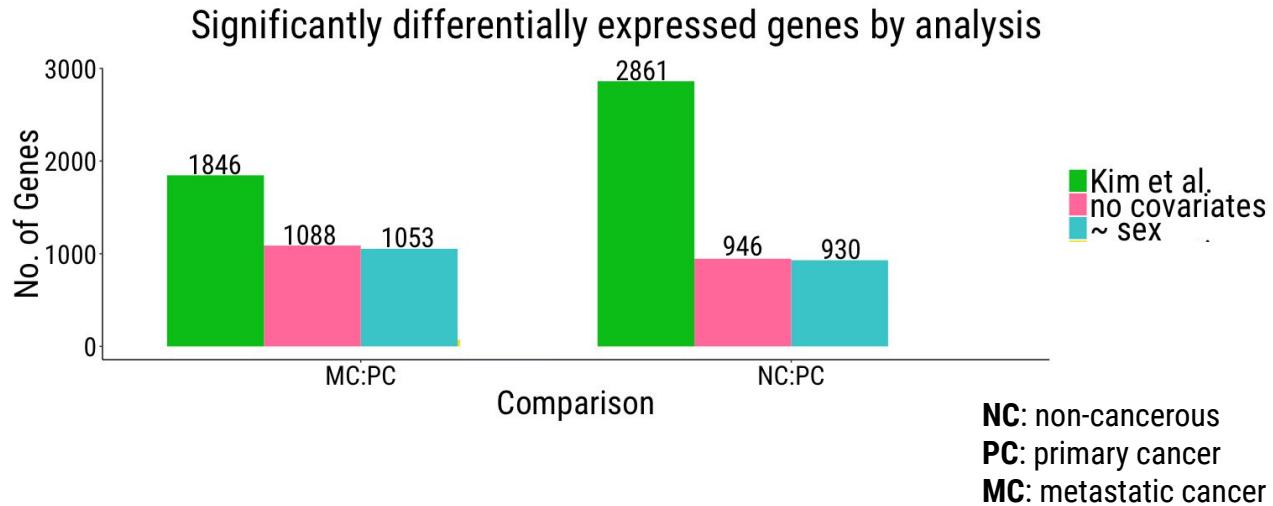
How similar are the results from Analysis A and Analysis B ?



Predictions can
be used to:

(1) Identify
studies of
interest

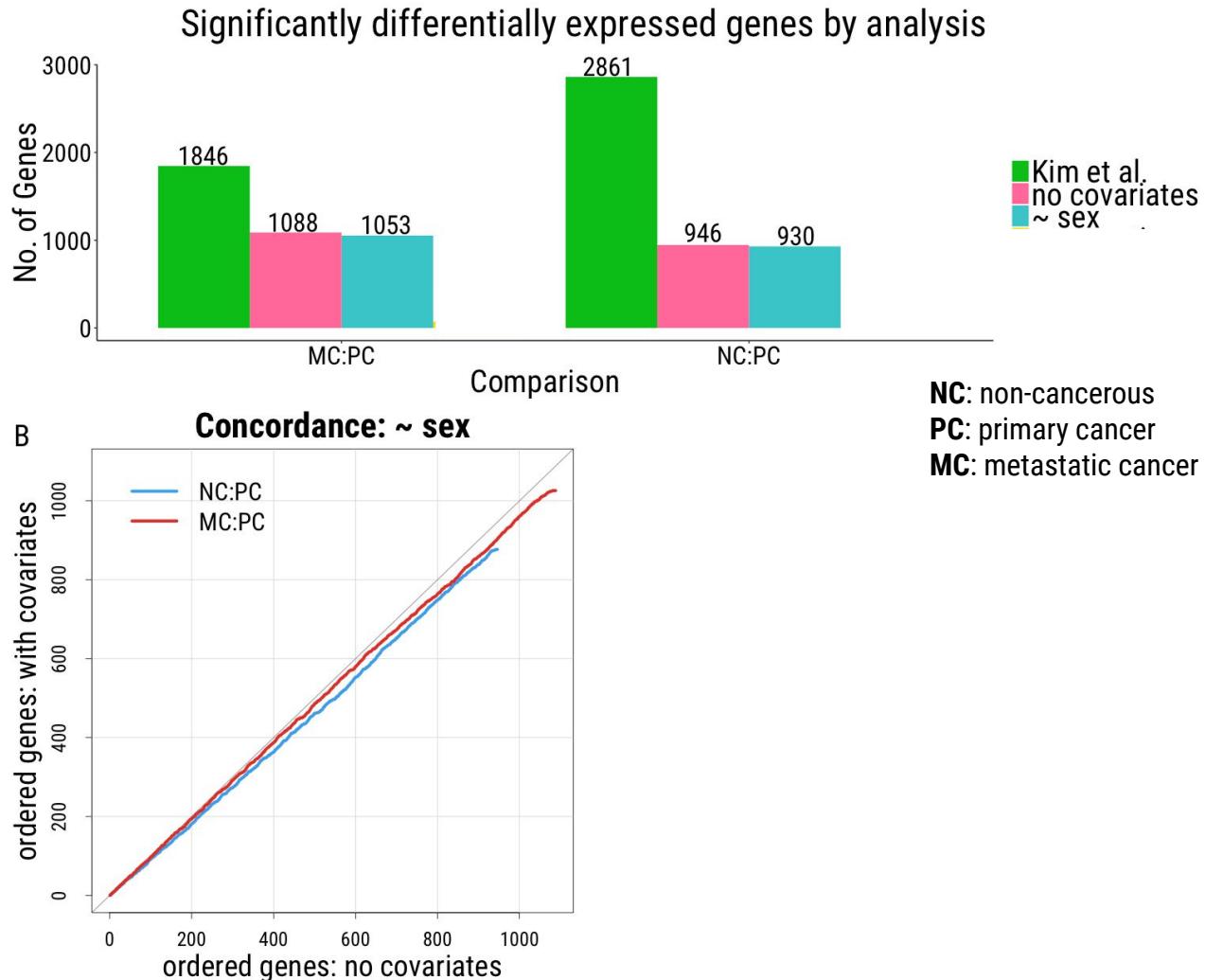
(2) appropriately
analyze data



Predictions can
be used to:

(1) Identify
studies of
interest

(2) appropriately
analyze data

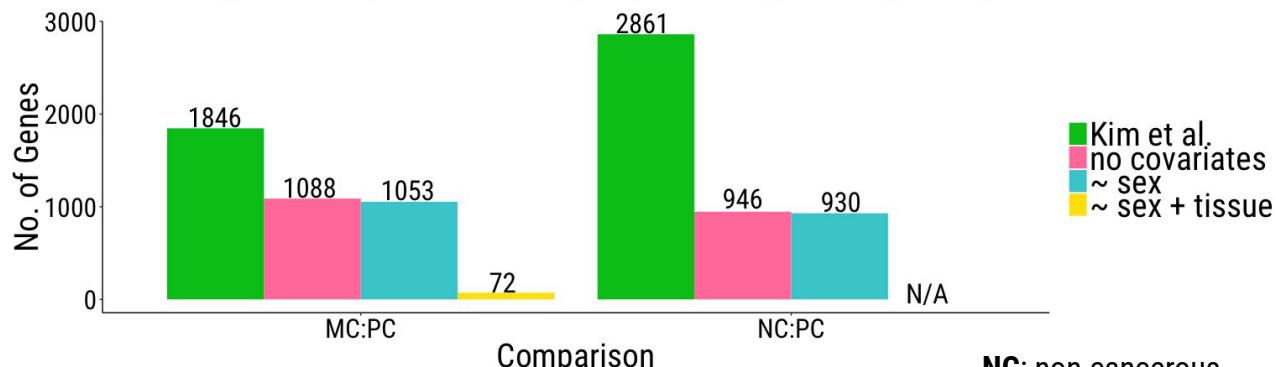


Predictions can
be used to:

(1) Identify
studies of
interest

(2) appropriately
analyze data

Significantly differentially expressed genes by analysis

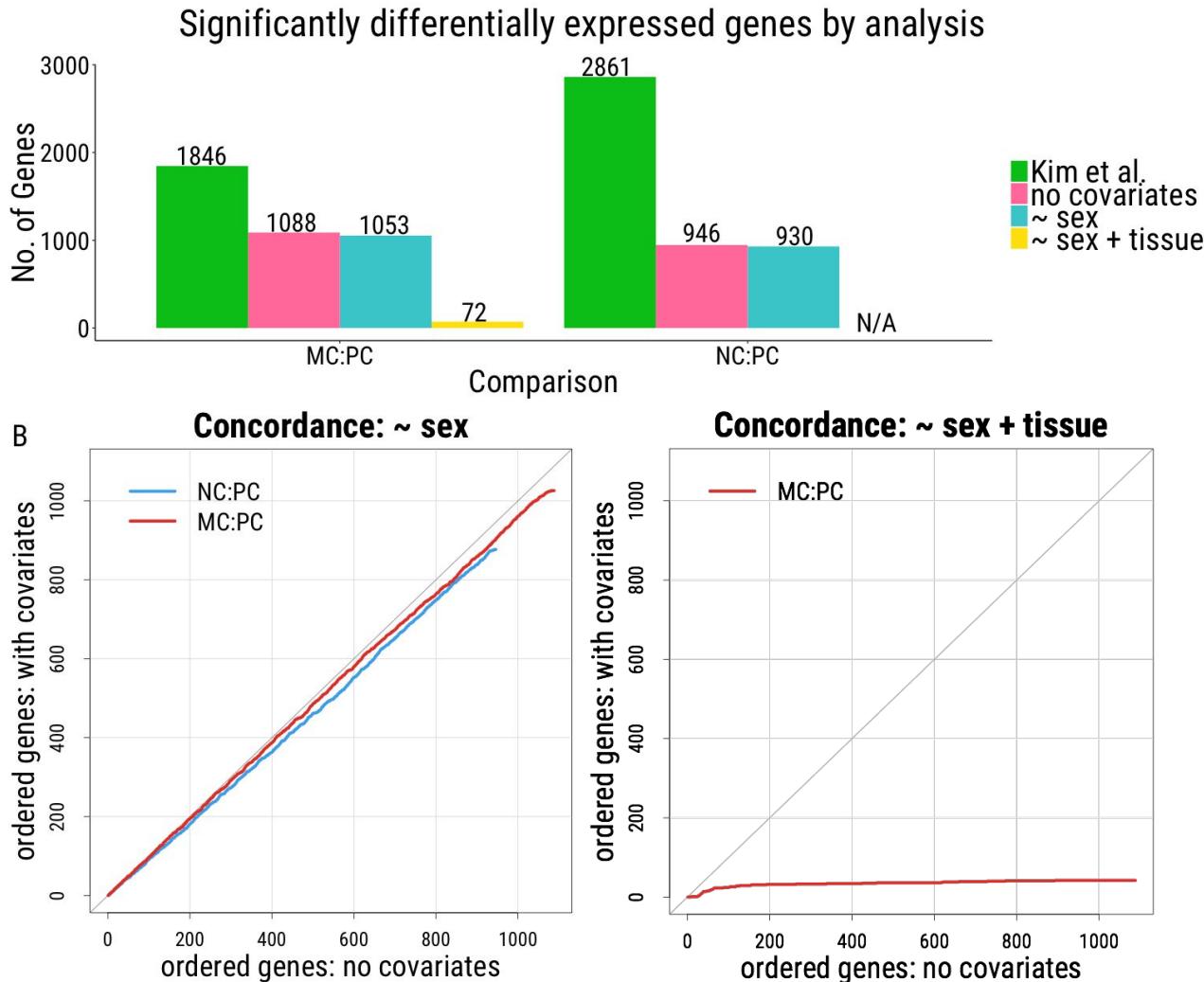


NC: non-cancerous
PC: primary cancer
MC: metastatic cancer

Predictions can
be used to:

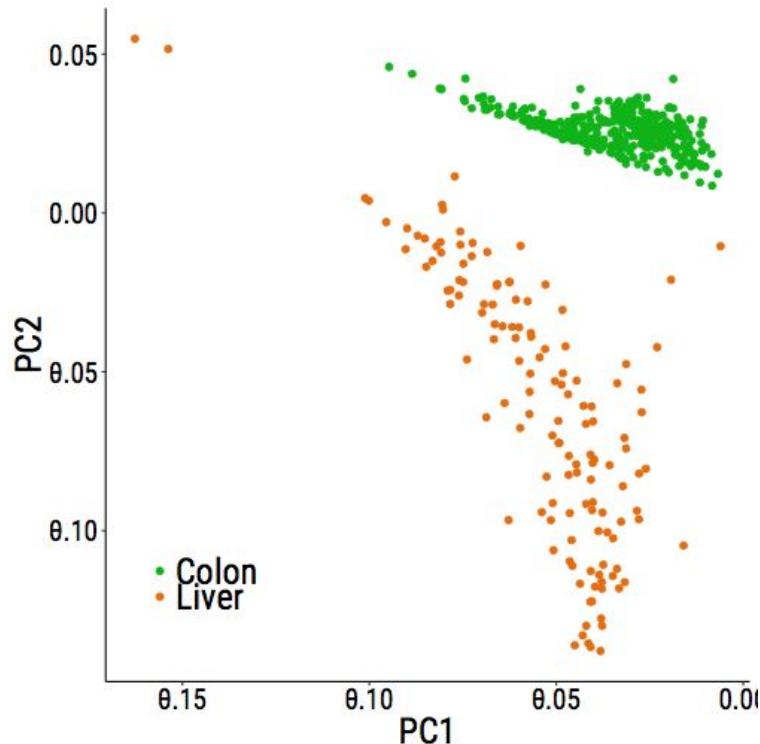
(1) Identify
studies of
interest

(2) appropriately
analyze data

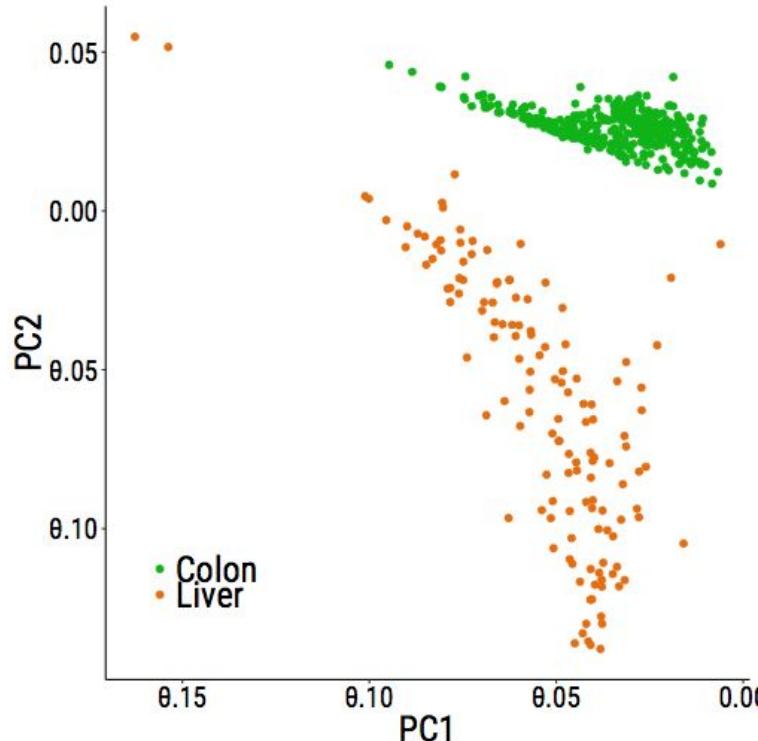


Loss of concordance suggests
that differential expression is
detecting tissue differences, not
cancer-related changes.

We have expression data from both healthy liver and colon samples (GTEx)...

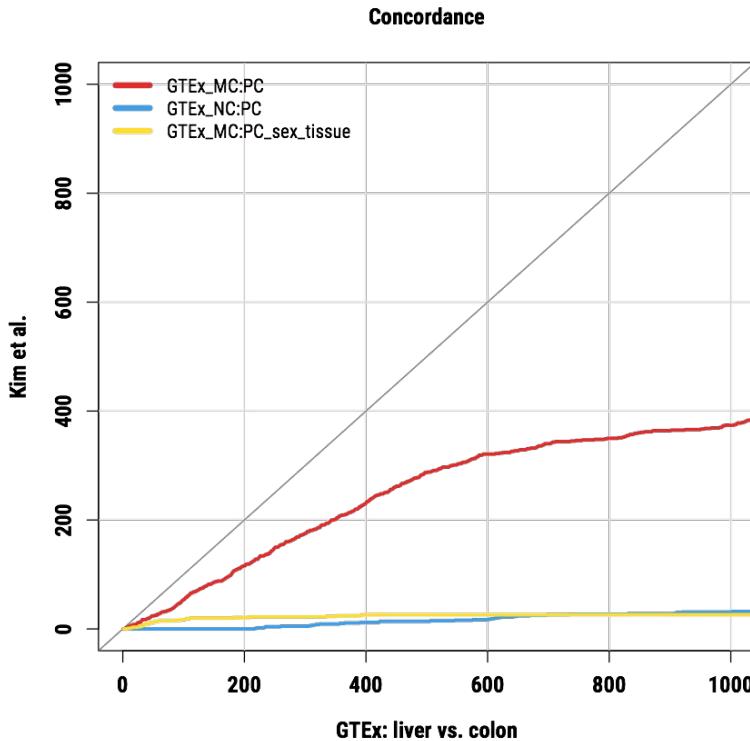
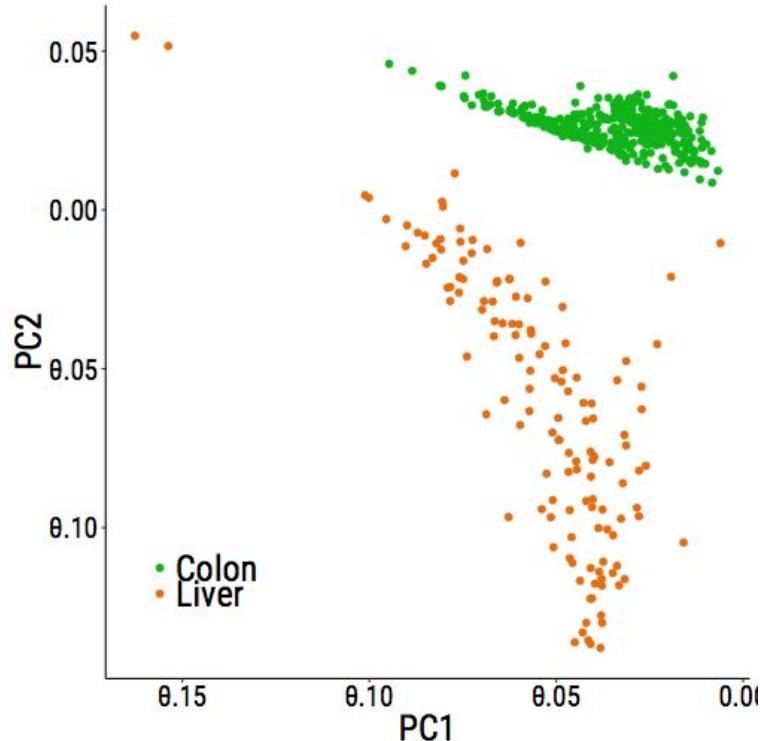


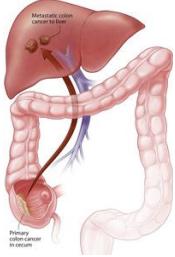
So...what if we compared the MC:PC results with differential expression between colon and liver?



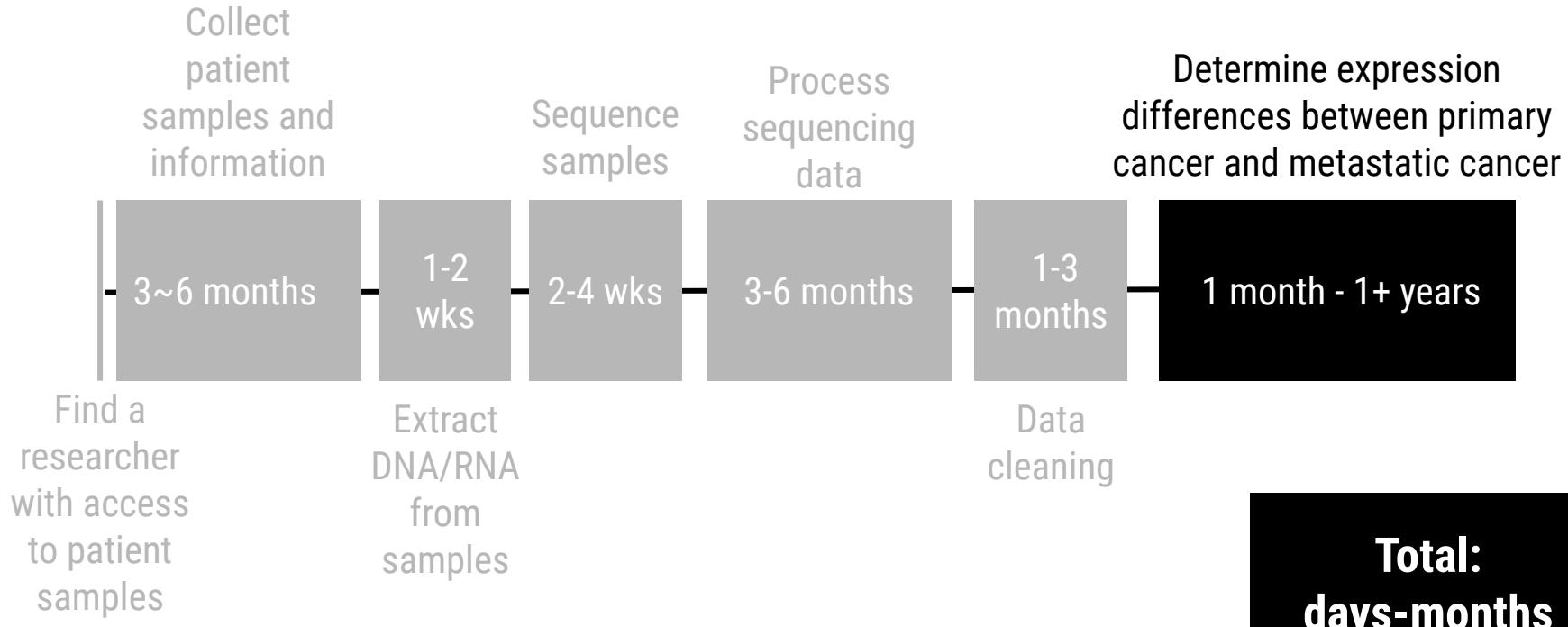
Hypothesis: MC:PC
results should be most
similar to GTEx colon
vs. liver

Comparison of results with GTEx colon vs. liver suggests differential expression results detecting tissue differences





What makes primary cancer different than metastatic cancer?



Finally, what if
YOU want to use
recount2...?

predictions (v0.0.06)

sample_id	dataset	reported_sex	predicted_sex	accuracy_sex	...	reported_tissue	predicted_tissue	accuracy_tissue
SRR660824	gtex	male	male	0.999	...	lung	lung	0.977
SRR2166176	gtex	male	male	0.999	...	brain	brain	0.977
SRR606939	gtex	female	female	0.999	...	heart	heart	0.966
SRR2167642	gtex	male	male	0.999	...	brain	brain	0.966
SRR2165473	gtex	male	male	0.999	...	skin	skin	0.966

Expression data and predictions available in recount R package

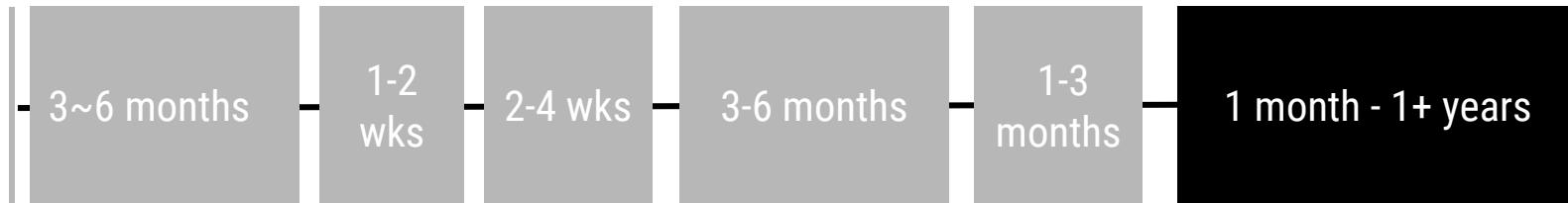
```
> library('recount')

> download_study('ERP001942', type='rse-gene')
> load(file.path('ERP001942', 'rse_gene.Rdata'))
> rse <- scale_counts(rse_gene)

> rse_with_pred <- add_predictions(rse)
```

There are *a lot* of questions an undergraduate could answer with *recount2*...

1. Which genes are expressed in which tissues?
2. Which genes contribute to X-Inactivation?
3. Which genes play a role in cancer? In autism? In Alzheimer's?
4. How different is expression between individuals?



Improving the value of public RNA-seq expression data by phenotype prediction

Shannon E Ellis, Leonardo Collado-Torres, Andrew Jaffe, Jeffrey T Leek 

Nucleic Acids Research, Volume 46, Issue 9, 18 May 2018, Page e54,

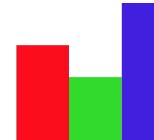
<https://doi.org/10.1093/nar/gky102>

Published: 05 March 2018 **Article history** ▾

<https://academic.oup.com/nar/article/46/9/e54/4920847>

...and all of this is published in Nucleic Acids Research

If you want to...

Align RNA-Seq data	 Rail-RNA Scalable RNA-seq alignment http://rail.bio
Learn about human expression	 recount2 https://jhbiostatistics.shinyapps.io/recount/
Predict phenotype information	 phenopredict https://github.com/leekgroup/phenopredict

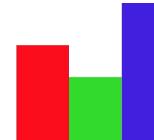
The Leek group

- Huan Chen
- Jack Fu
- Aboozar Hadavand
- Leslie Myint
- Kayode Sosina
- Sara Wang
- **Jeff Leek**

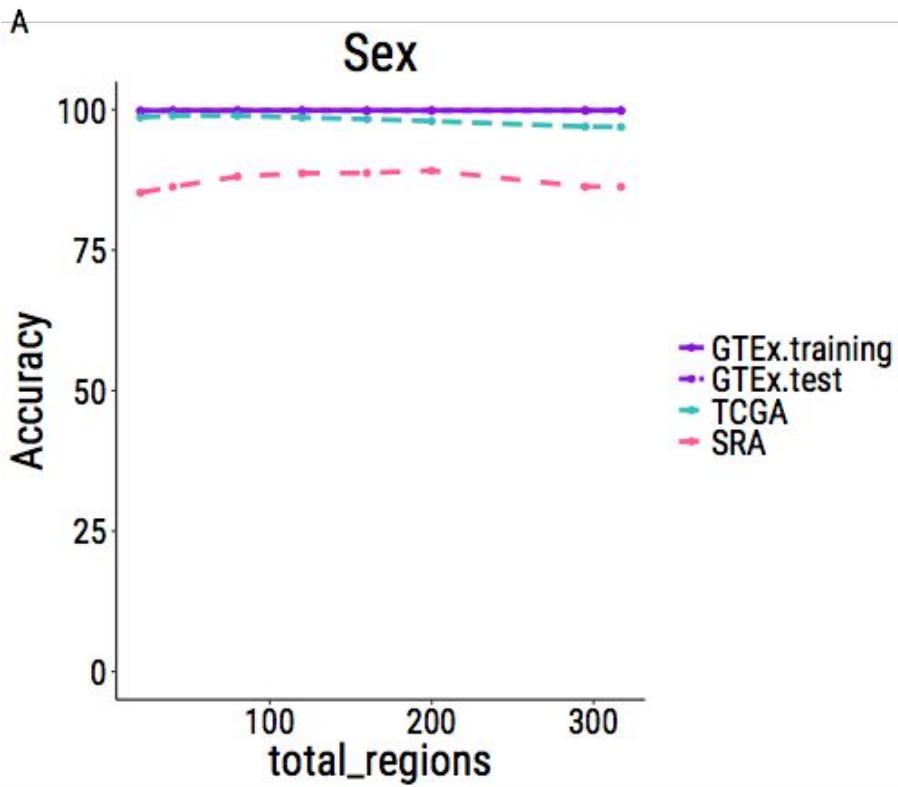
Collaborators

- Andrew Jaffe
- Kasper Hansen
- Margaret Taub
- Leah Jager
- Ben Langmead
- Abhi Nellore
- Kai Kammers
- Leo Collado-Torres
- Ashkaun Razmara

If you want to...

Align RNA-Seq data	 Rail-RNA Scalable RNA-seq alignment <i>http://rail.bio</i>
Learn about human expression	 recount2 <i>https://jhbiostatistics.shinyapps.io/recount/</i>
Predict phenotype information	 phenopredict <i>https://github.com/leekgroup/phenopredict</i>

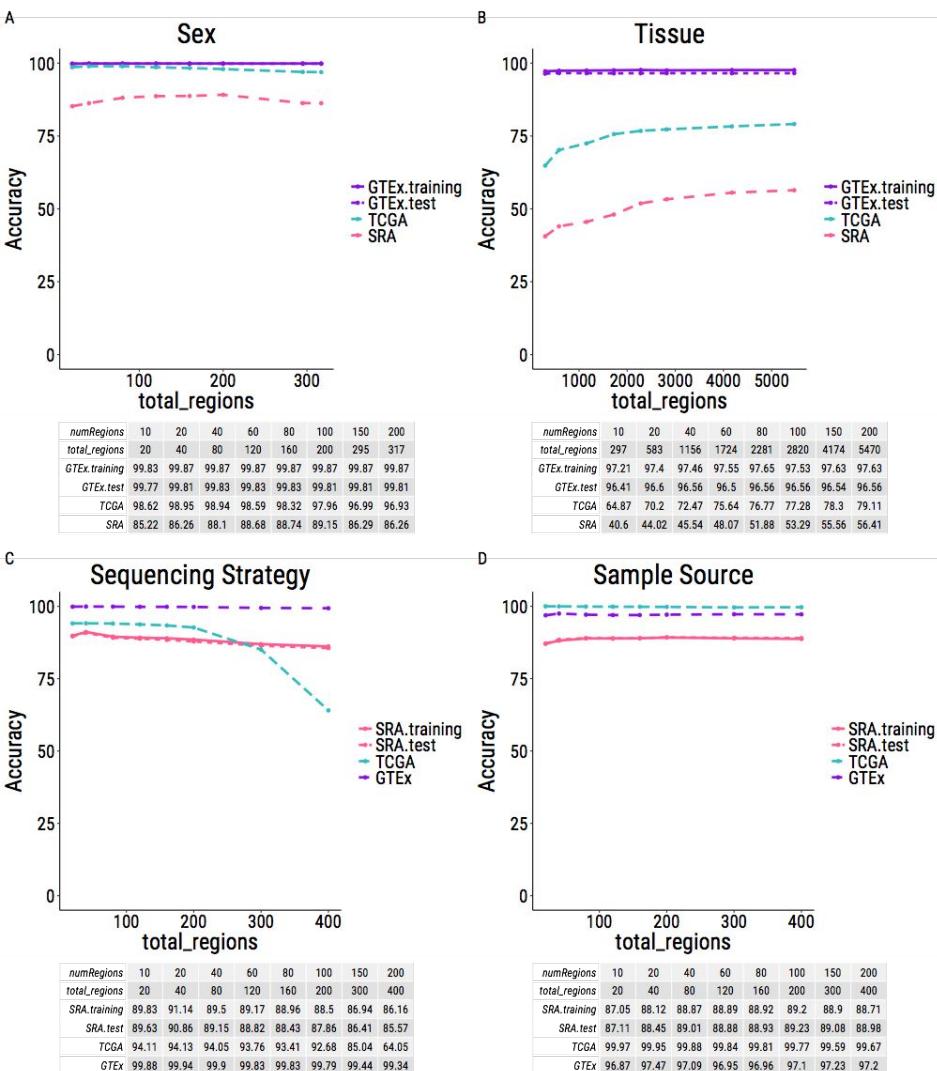
Number of regions used for prediction was optimized in the training set.



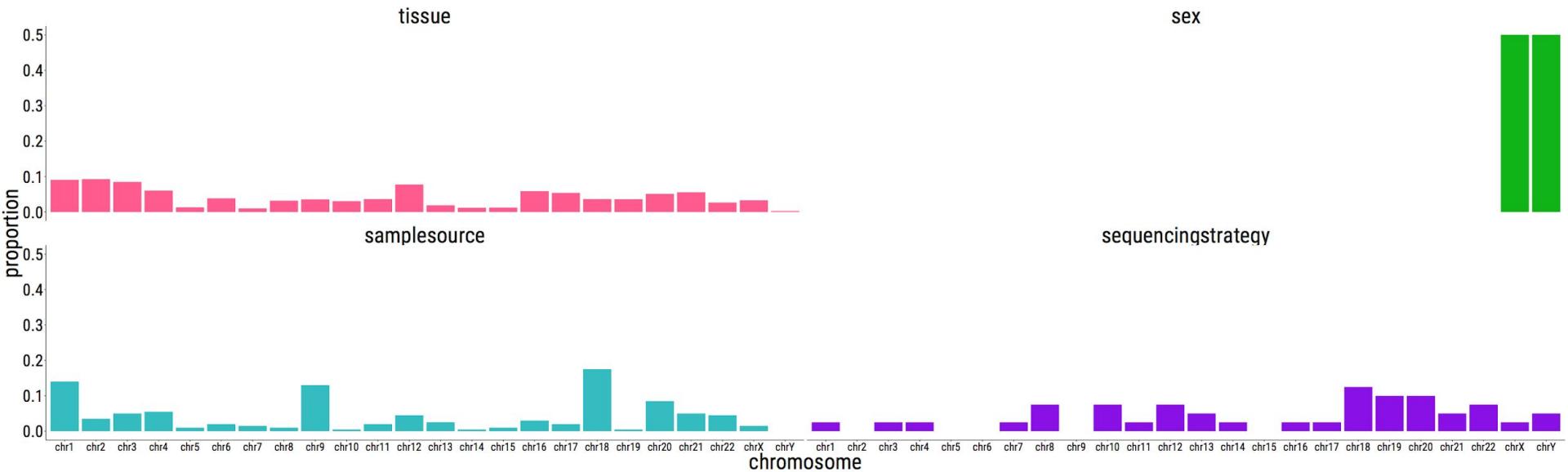
numRegions	10	20	40	60	80	100	150	200
total_regions	20	40	80	120	160	200	295	317
GTEx.training	99.83	99.87	99.87	99.87	99.87	99.87	99.87	99.87
GTEx.test	99.77	99.81	99.83	99.83	99.83	99.81	99.81	99.81
TCGA	98.62	98.95	98.94	98.59	98.32	97.96	96.99	96.93
SRA	85.22	86.26	88.1	88.68	88.74	89.15	86.29	86.26

C

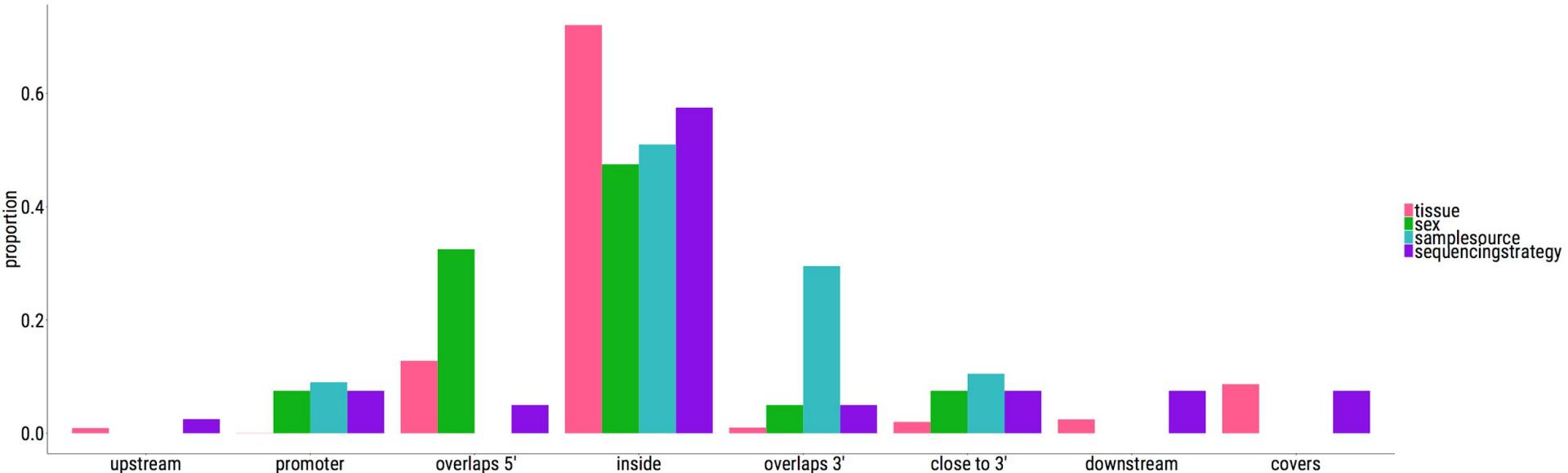
Phenotypes are largely insensitive to number of regions used to build predictor

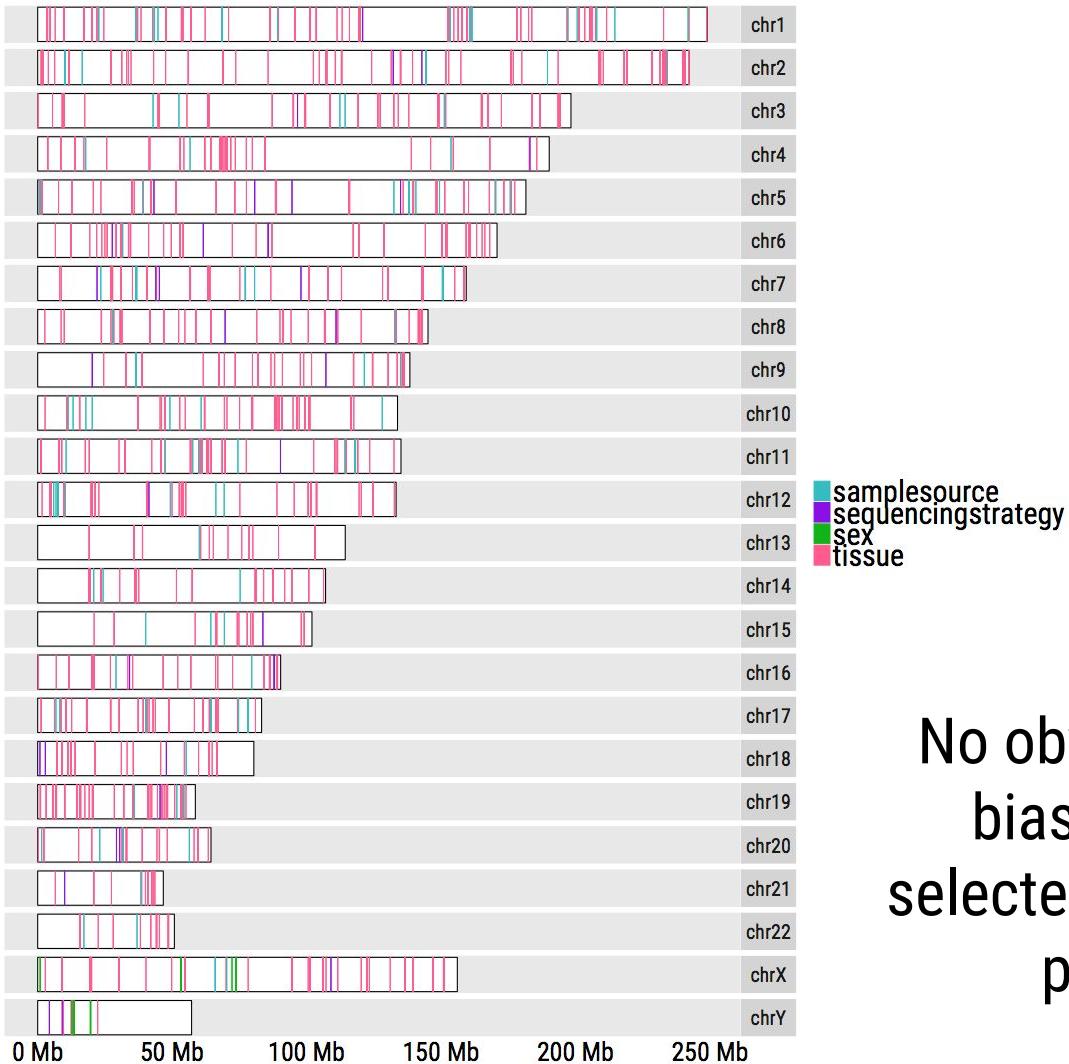


Aside from sex prediction, regions were selected across all chromosomes



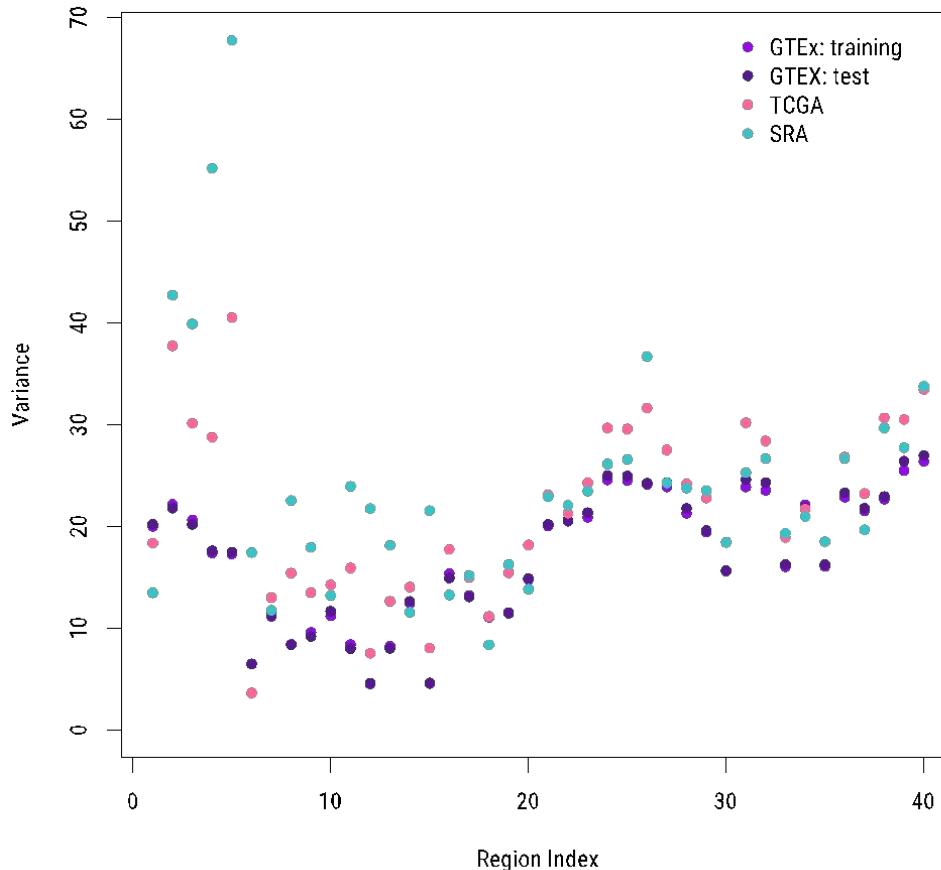
Most regions used for prediction fell within the annotated gene; however sequencing strategy shows some bias for regions outside the annotated transcriptome.



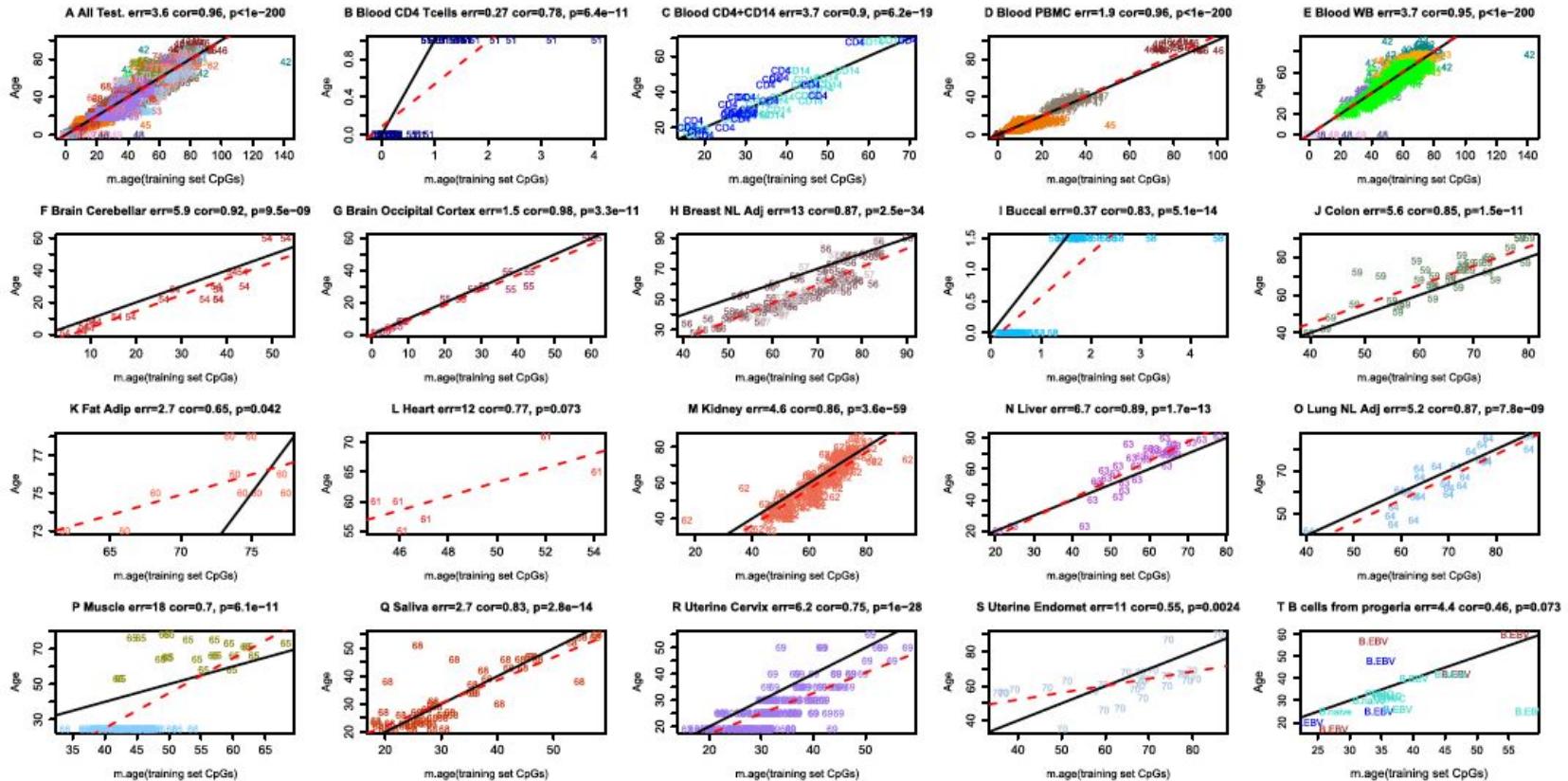


No obvious location bias for regions selected for any of the predictors

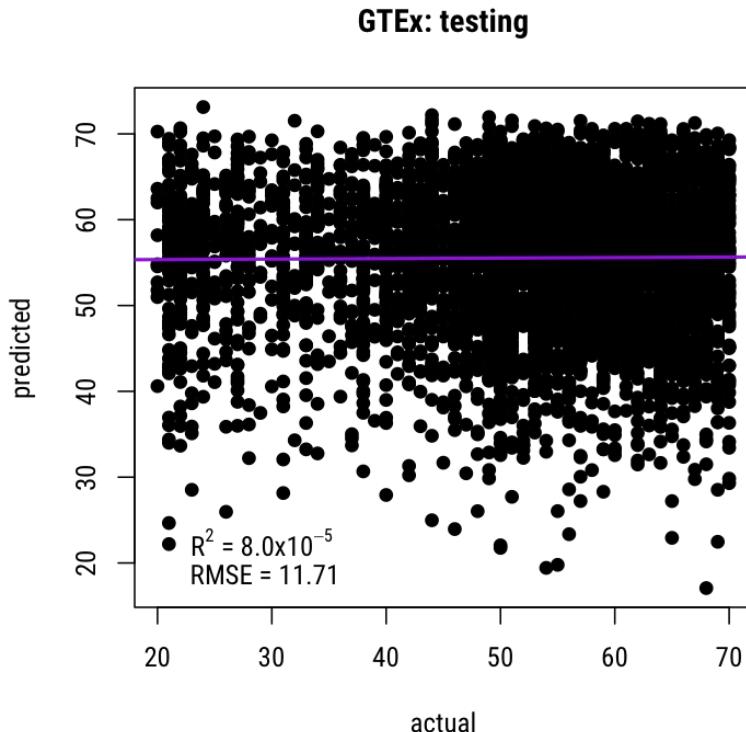
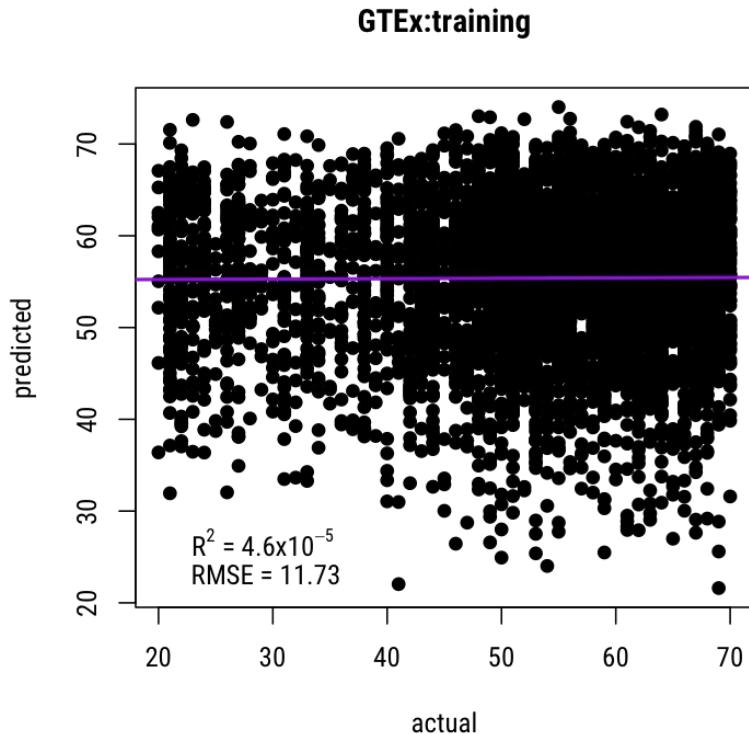
Variance across sex predicting regions



Horvath demonstrates that 353 CpGs can accurately predict age

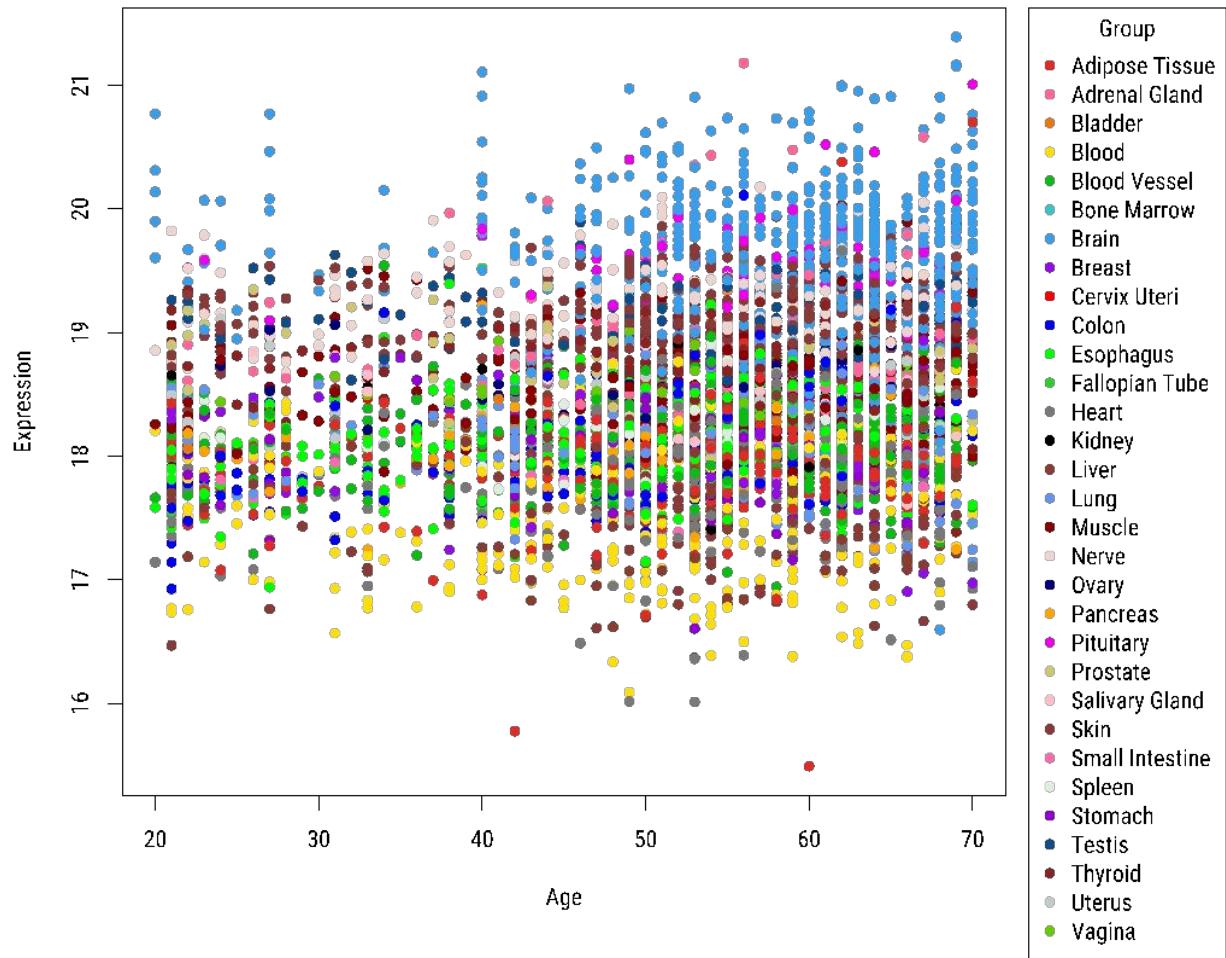


How well can we predict age in GTEx?

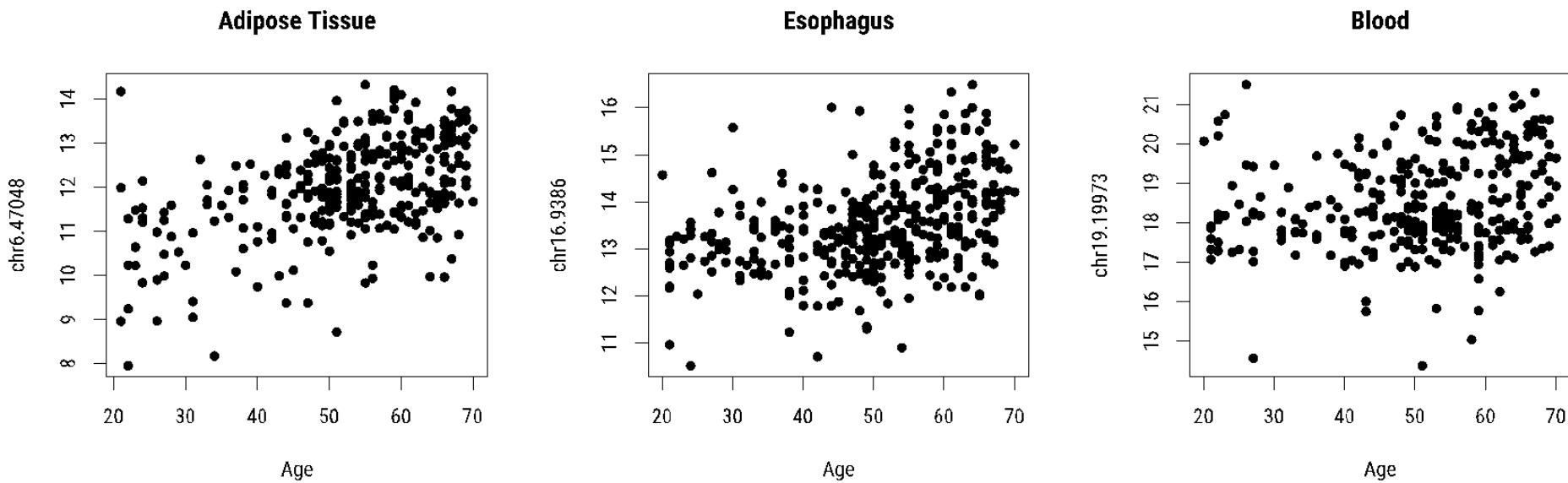


Tissue poses a problem for prediction...

chr16.2196

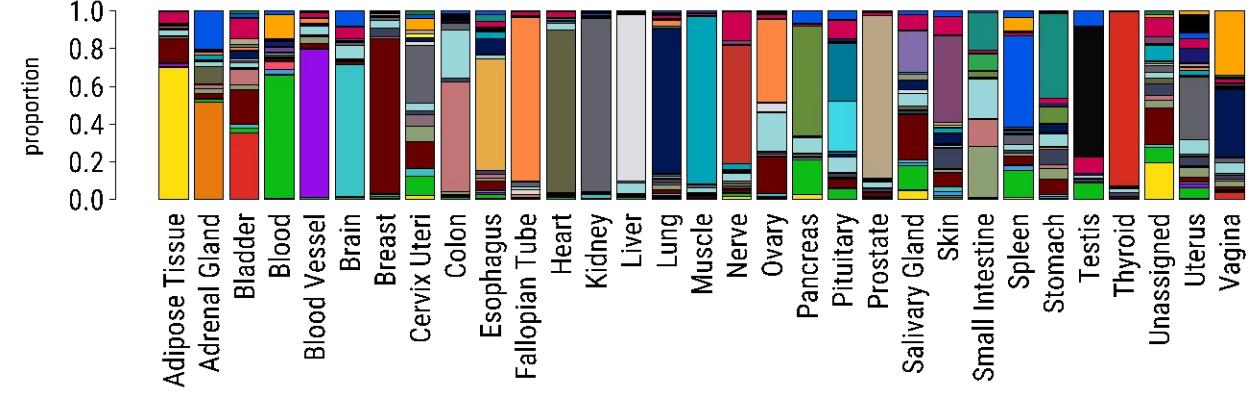


Even within tissue, signal is pretty weak...



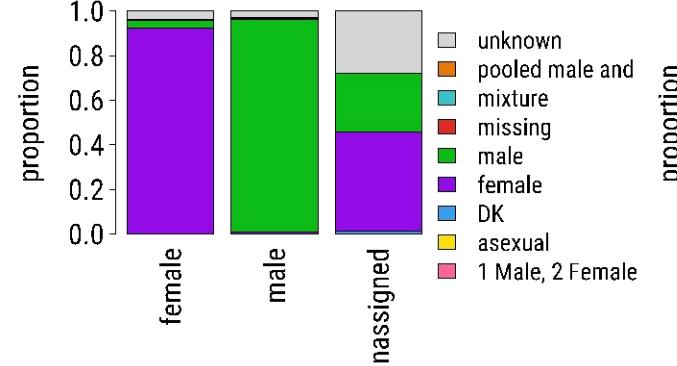
A

Tissue



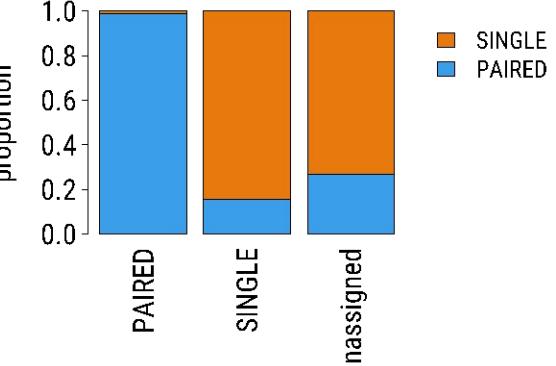
B

Sex



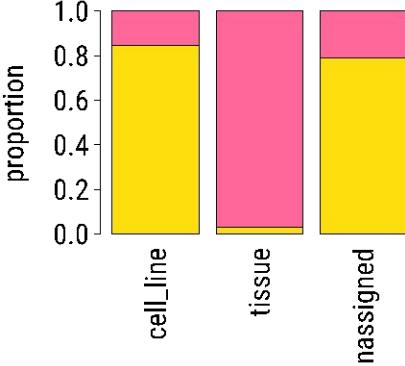
C

Sequencing Strategy

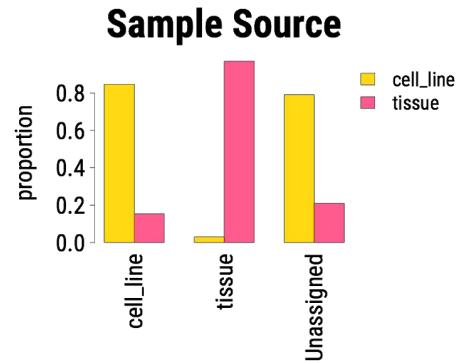
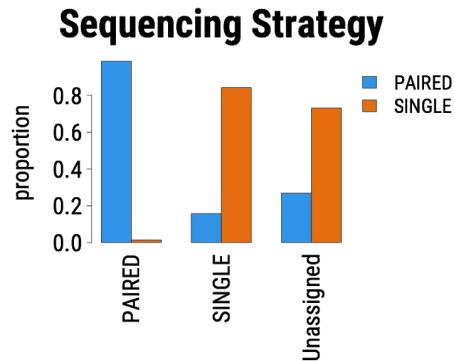
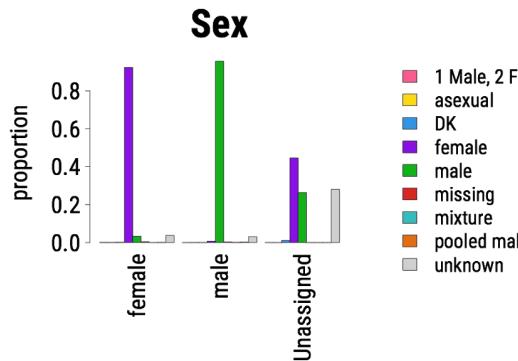
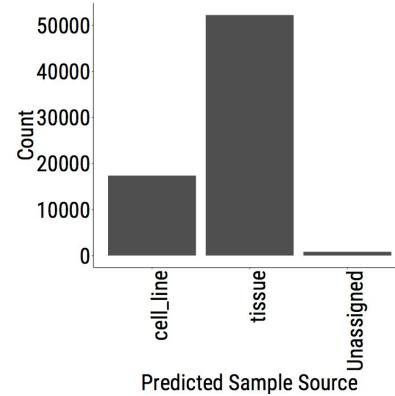
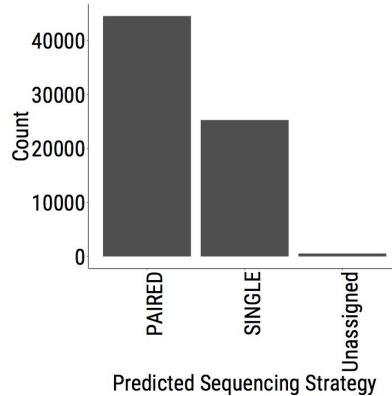
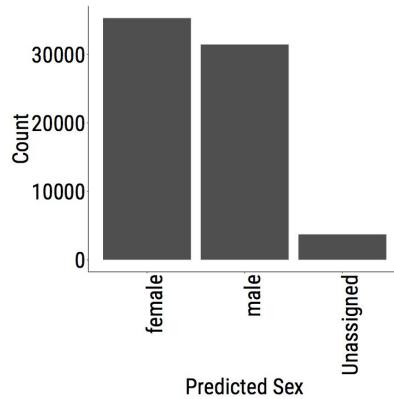


D

Sample Source



When we look at the other three predictions...



When we look at the other three predictions...

