

# Working with a whole bunch of genetic data ...and teaching data science

Shannon E. Ellis  
Assistant Teaching Professor  
COGS1 (Spring 2019)

# A quick tour of a geneticist turned data scientist

## Background

## Projects

1. PhD work studying the genetic basis of autism
2. Postdoctoral work working with 70,000 samples
3. Working toward accessible data science education

## What I do here at UCSD

SHANNON



I LIKE ALL THE THINGS!

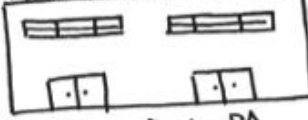


SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA



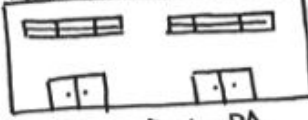
SCIENCE IS PRETTY COOL!

SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

KING'S COLLEGE



Wilkes-Barre, PA



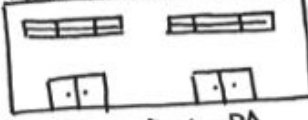
SCIENCE IS PRETTY COOL!

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!

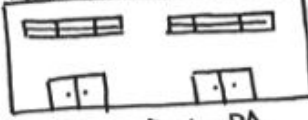


SHANNON



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!

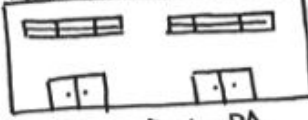


SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

GENETICS IS AWESOME!

...but what did that software actually do?



KING'S COLLEGE



Wilkes-Barre, PA

SCIENCE IS PRETTY COOL!



JOHNS HOPKINS



Baltimore, MD

DATA ANALYSIS IS WHERE IT'S AT!



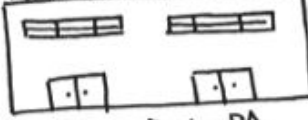


SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

SCIENCE IS PRETTY COOL!



KING'S COLLEGE



Wilkes-Barre, PA

GENETICS IS AWESOME!

...but what did that software actually do?



JOHNS HOPKINS



Baltimore, MD

DATA ANALYSIS IS WHERE IT'S AT!



JEFF LEEK



COME DO A POSTDOC WITH ME!

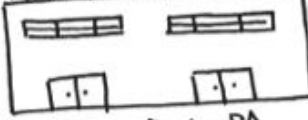


SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

SCIENCE IS PRETTY COOL!



KING'S COLLEGE



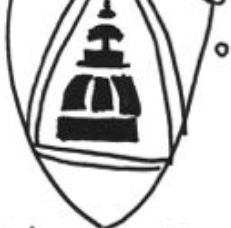
Wilkes-Barre, PA

GENETICS IS AWESOME!

...but what did that software actually do!



JOHNS HOPKINS



Baltimore, MD

I LOVE TEACHING!



DATA ANALYSIS IS WHERE IT'S AT!



JEFF LEEK



COME DO A POSTDOC WITH ME!

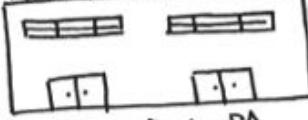


SHANNON



I LIKE ALL THE THINGS!

HIGH SCHOOL O' MANY



Upper Darby, PA

SCIENCE IS PRETTY COOL!



KING'S COLLEGE



Wilkes-Barre, PA

GENETICS IS AWESOME!

...but what did that software actually do?



JOHNS HOPKINS



Baltimore, MD

I LOVE TEACHING!



DATA ANALYSIS IS WHERE IT'S AT!



JEFF LEEK



COME DO A POSTDOC WITH ME!



TODAY

AN OPEN SCIENCE



# A quick tour of a geneticist turned data scientist

## Background

## Projects

1. PhD work studying the genetic basis of autism
2. Postdoctoral work working with 70,000 samples
3. Working toward accessible data science education

## What I do here at UCSD

# The quickest history of human genetics

1900 - Rediscovery of Mendel's work

1950s - DNA is the genetic material & Structure of DNA

2001 - Human Genome Project

- Thought this would unlock understanding of all disease!
- Not exactly how things worked out - way more complex than initially thought

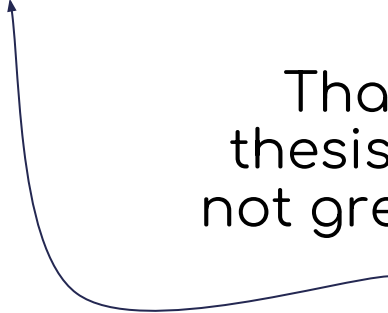
Mid-2000s - GWAS set out to sort all this out

- 2007: 240 papers published
- 2017: 3800+

2010s - a move toward additional approaches

# Multi-omic Data Provide a More Complete Understanding of the Autistic Brain

That's the title of my thesis dissertation. I'm not great at coming up with titles...ever.



# General Outline

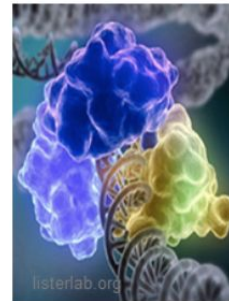
## I. Autism Background

## II. Transcriptome Analyses

- A. Gene expression differences in the autistic brain
- B. Cross-disorder transcriptomic overlap

## III. Epigenome of the Autistic Brain

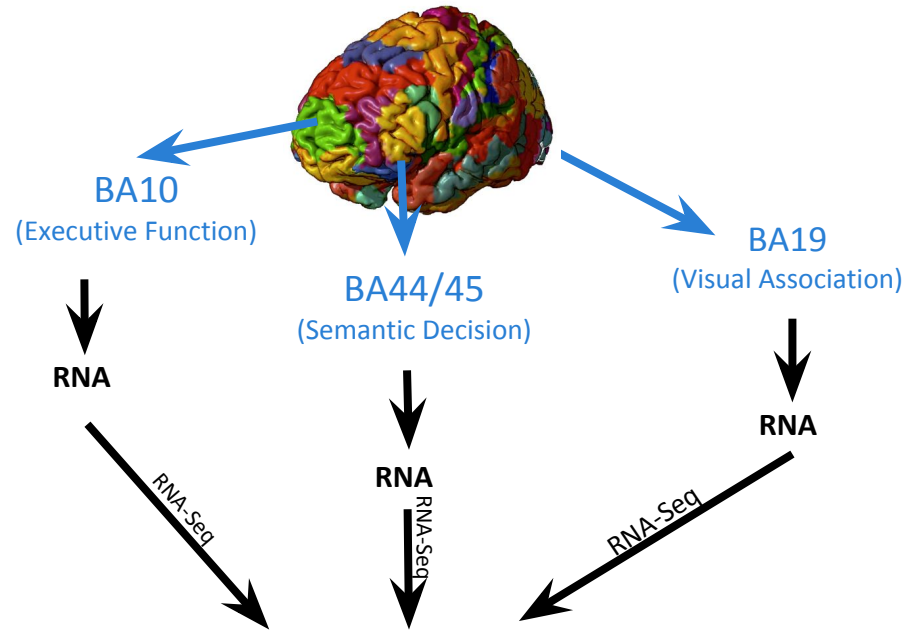
- A. CpG methylation
- B. nonCpG methylation



# RNA-Sequencing in Autism Brains



Simone Gupta, PhD



## ARTICLE

Received 28 Sep 2014 | Accepted 3 Nov 2014 | Published 10 Dec 2014

DOI: 10.1038/ncomms6748

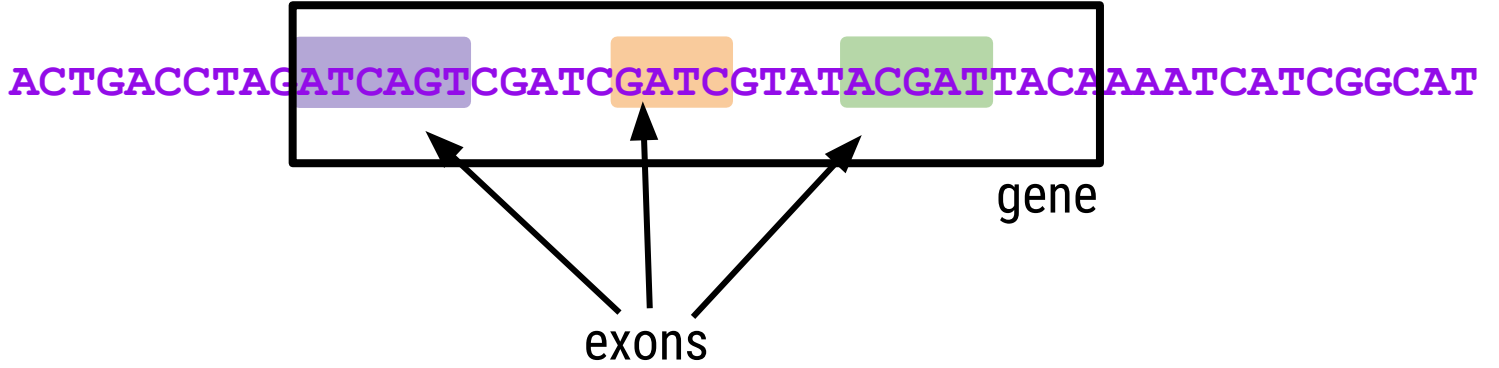
OPEN

Transcriptome analysis reveals dysregulation of innate immune response genes and neuronal activity-dependent genes in autism

Simone Gupta<sup>1</sup>, Shannon E. Ellis<sup>1</sup>, Foram N. Ashar<sup>1</sup>, Anna Moes<sup>1</sup>, Joel S. Bader<sup>1,2</sup>, Jianan Zhan<sup>2</sup>, Andrew B. West<sup>3</sup> & Dan E. Arking<sup>1</sup>



# The Central Dogma of Genetics



# The Central Dogma of Genetics

DNA



ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT



**transcription**

RNA



AUCAGUCGAUCACCGAU

# The Central Dogma of Genetics

**DNA**



ACTGACCTAGATCAGTCGATCGATCGTATACGATTACAAAATCATCGGCAT



**transcription**

**RNA**

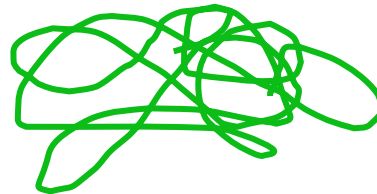


AUCAGUCGAUCACCGAU




**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>
<b>DNA</b> 	<i>blueprint</i>	<i>2</i>	<i>gene</i>	<i>20,000</i>



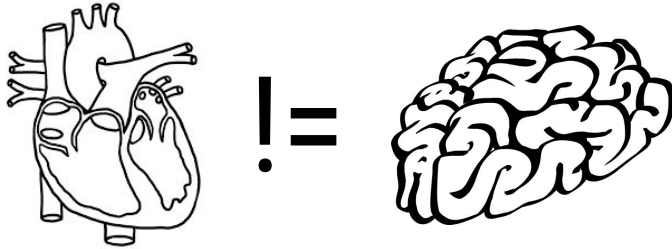
# Two copies of **DNA** -> many **transcripts** -> many **proteins**

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

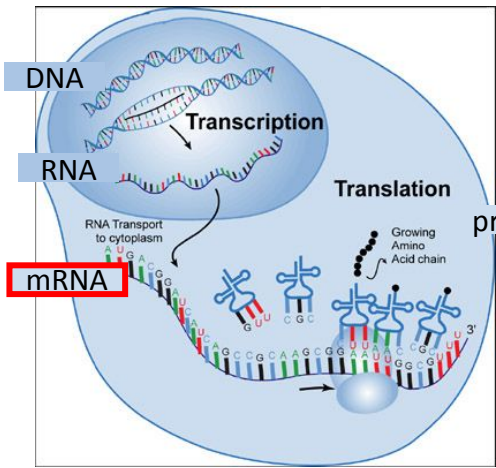
# Two copies of **DNA** -> many **transcripts** -> many **proteins**

	<i>role in the cell</i>	<i># copies/cell</i>	<i>functional unit</i>	<i># unique functional units</i>
<b>DNA</b> 	<i>blueprint</i>	<i>2</i>	<i>gene</i>	<i>20,000</i>
<b>RNA</b> 	<i>messenger</i>	<i>varies ~360,000</i>	<i>transcript</i>	<i>~100,000</i>
<b>proteins</b> 	<i>carry out cellular functions</i>	<i>varies ~10<sup>10</sup></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

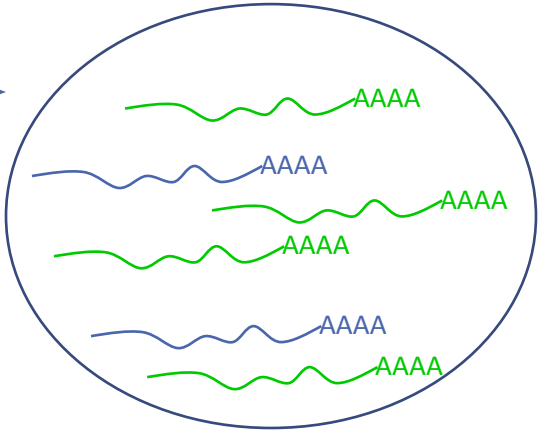


# RNA-Sequencing



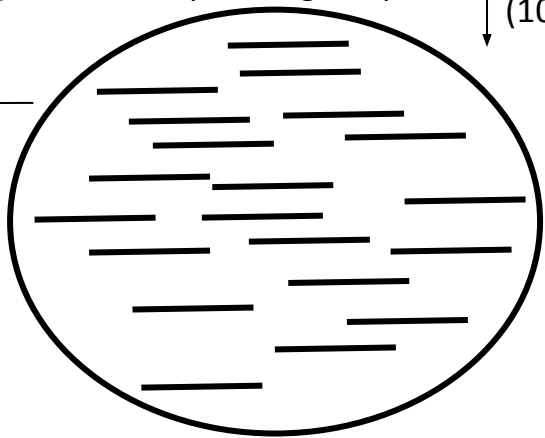
Extract mRNA

Pool of polyA+ -selected mRNA

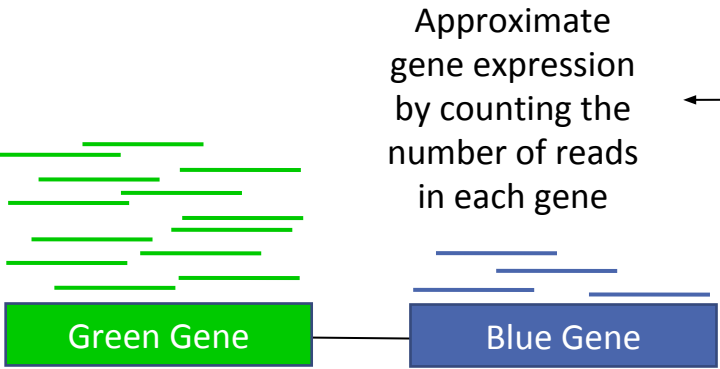


Prepare a sequencing library

Illumina Sequencing (100bp SE)



Align those reads to the reference genome



Green Gene > Blue Gene



# Results in a single slide

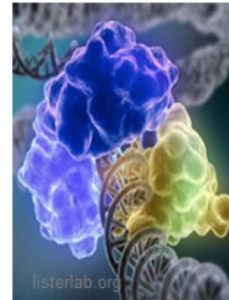
## I. Autism Background

## II. Transcriptome Analyses

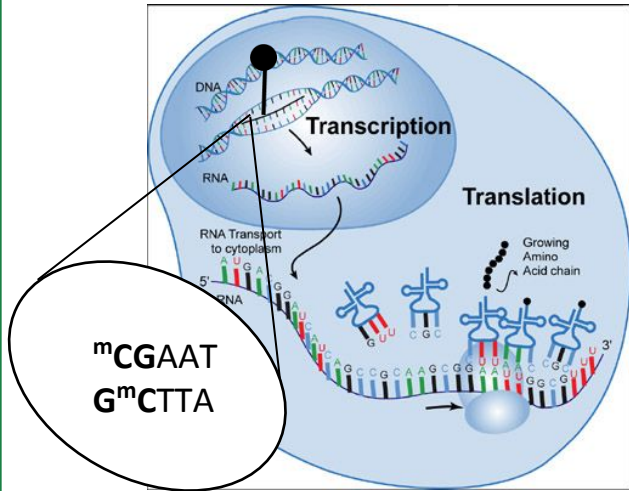
- A. Microglia playing a role in the autistic brain
- B. RNA levels show similar patterns across conditions

## III. Epigenome of the Autistic Brain

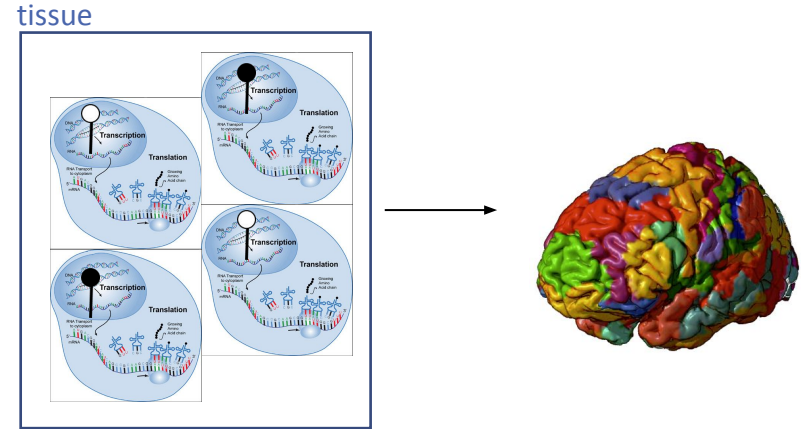
- A. CpG methylation
- B. nonCpG methylation



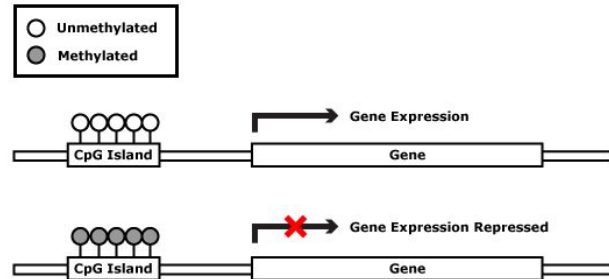
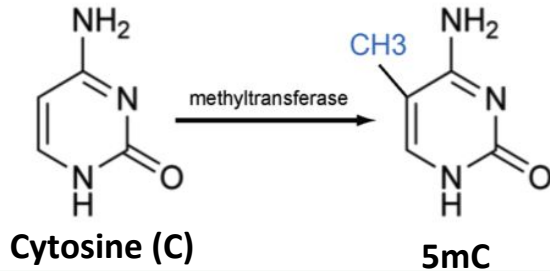
# DNA methylation is most often studied at CpG dinucleotides



<http://www.tokresource.org/>



50% methylated



# Results in a single slide

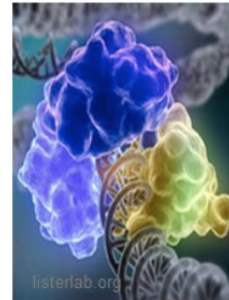
## I. Autism Background

## II. Transcriptome Analyses

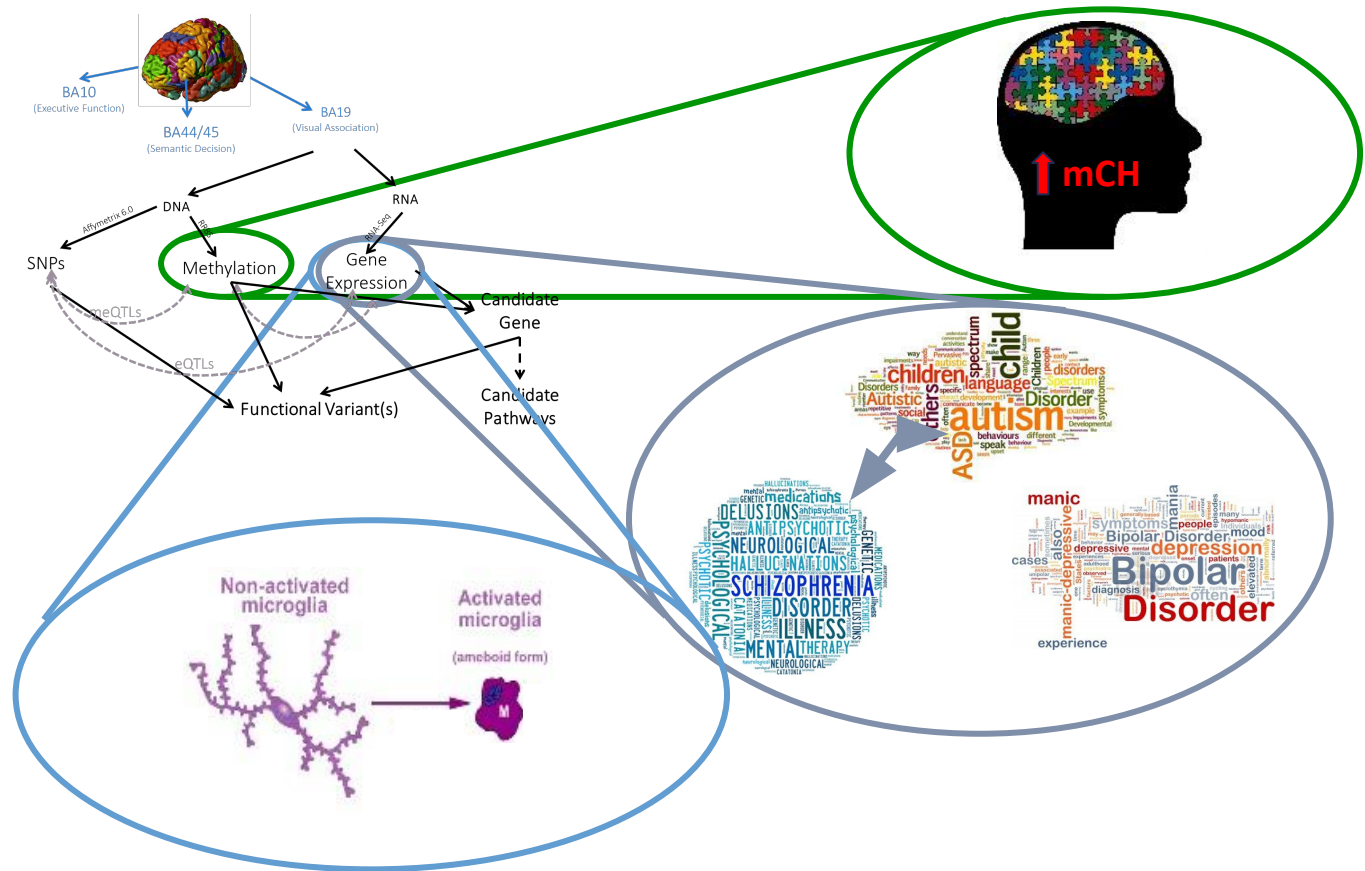
- A. Microglia playing a role in the autistic brain
- B. RNA levels show similar patterns across conditions

## III. Epigenome of the Autistic Brain

- A. CpG methylation does not differ
- B. Increased global nonCpG methylation



# Conclusions: Toward a More Complete Understanding of the Autistic Brain



# Scientific Acknowledgments (PhD work)

## The Arking Lab



Dan E. Arking, PhD



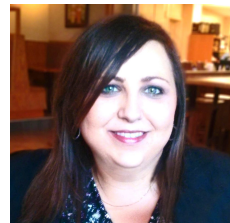
Foram N. Ashar



Nathan Bilhmayer



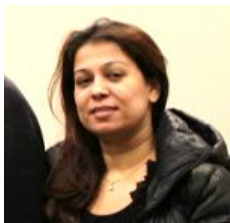
Pallav Bhatnagar, PhD



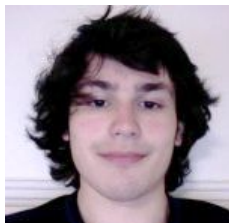
Christina Castellani, PhD



Rebecca Eggebeen



Simone Gupta, PhD



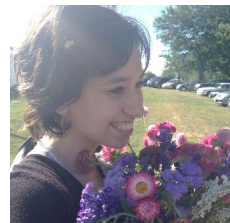
Naftali Horowitz



Ryan Longchamps



Anna Moes, MS



Rebecca Panitch



Elizabeth Vincent

## Collaborators



Shan Andrews

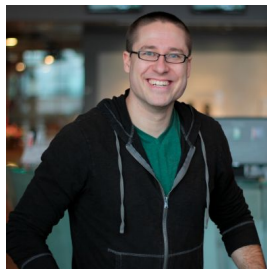


Andrew West, PhD

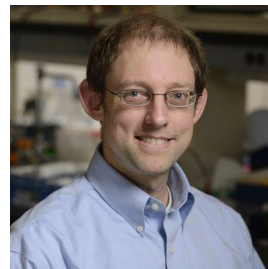
## Thesis Committee



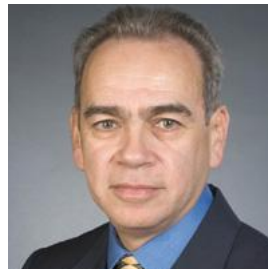
Dani Fallin, PhD



Jeff Leek, PhD



Joel Bader, PhD



Juan Troncoso, MD

## Human Genetics Program



David Valle, MD



Kirby Smith, PhD



Sandy Muscelli

# A quick tour of a geneticist turned data scientist

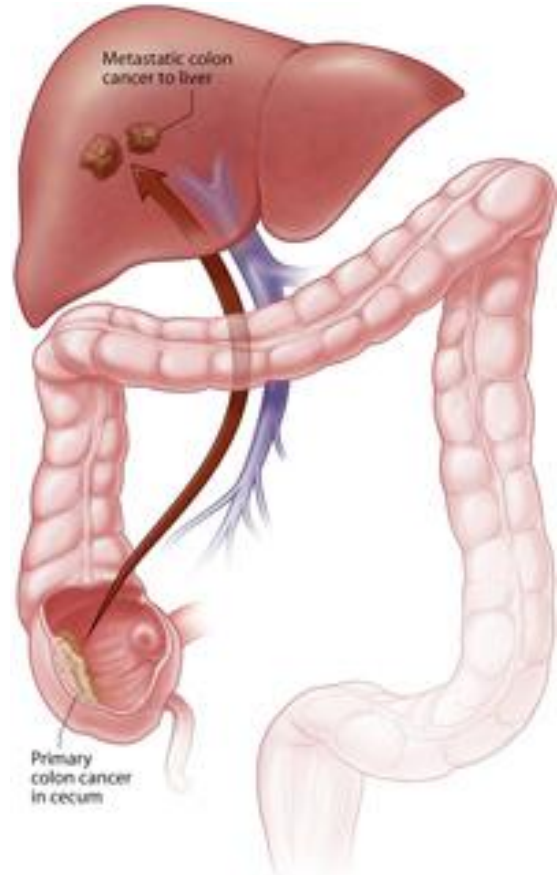
## Background

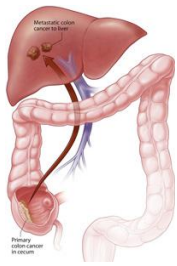
## Projects

1. PhD work studying the genetic basis of autism
2. Postdoctoral work working with 70,000 samples
3. Working toward accessible data science education

## What I do here at UCSD

What makes primary cancer different than metastatic cancer?



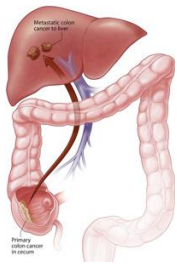


# What makes primary cancer different than metastatic cancer?

---

Find a  
researcher  
with access  
to patient  
samples



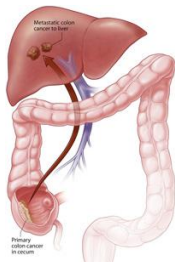


# What makes primary cancer different than metastatic cancer?

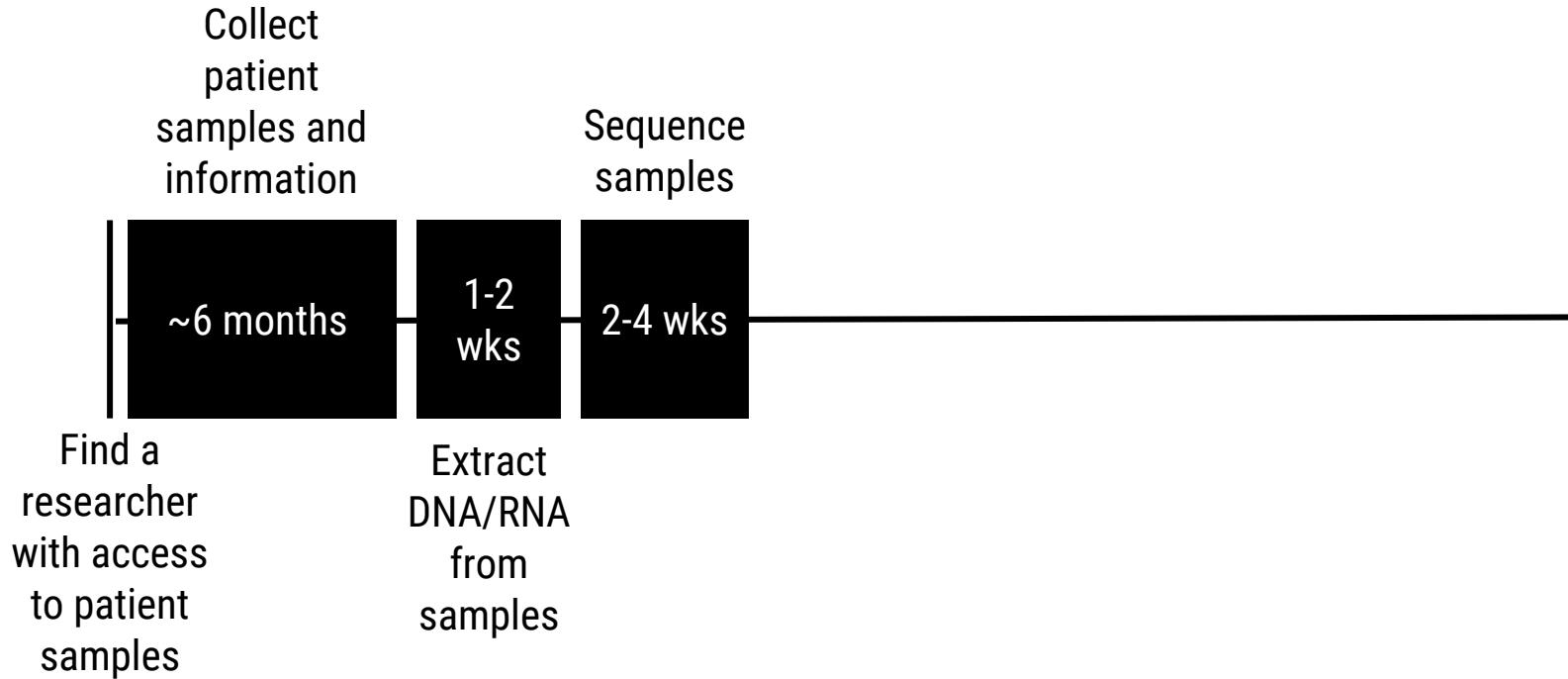
Collect  
patient  
samples and  
information

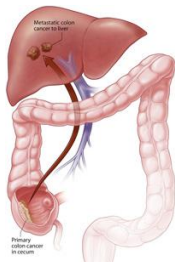
~6 months

Find a  
researcher  
with access  
to patient  
samples

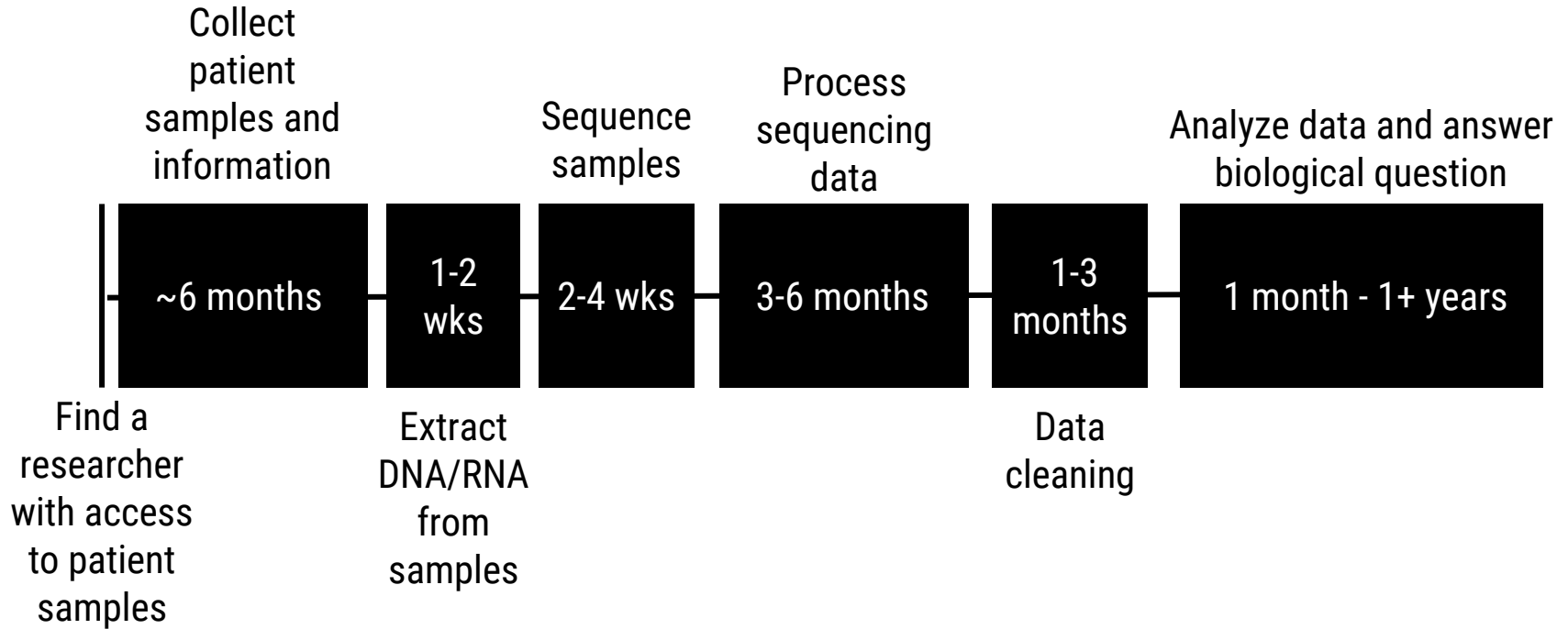


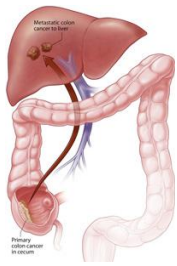
# What makes primary cancer different than metastatic cancer?



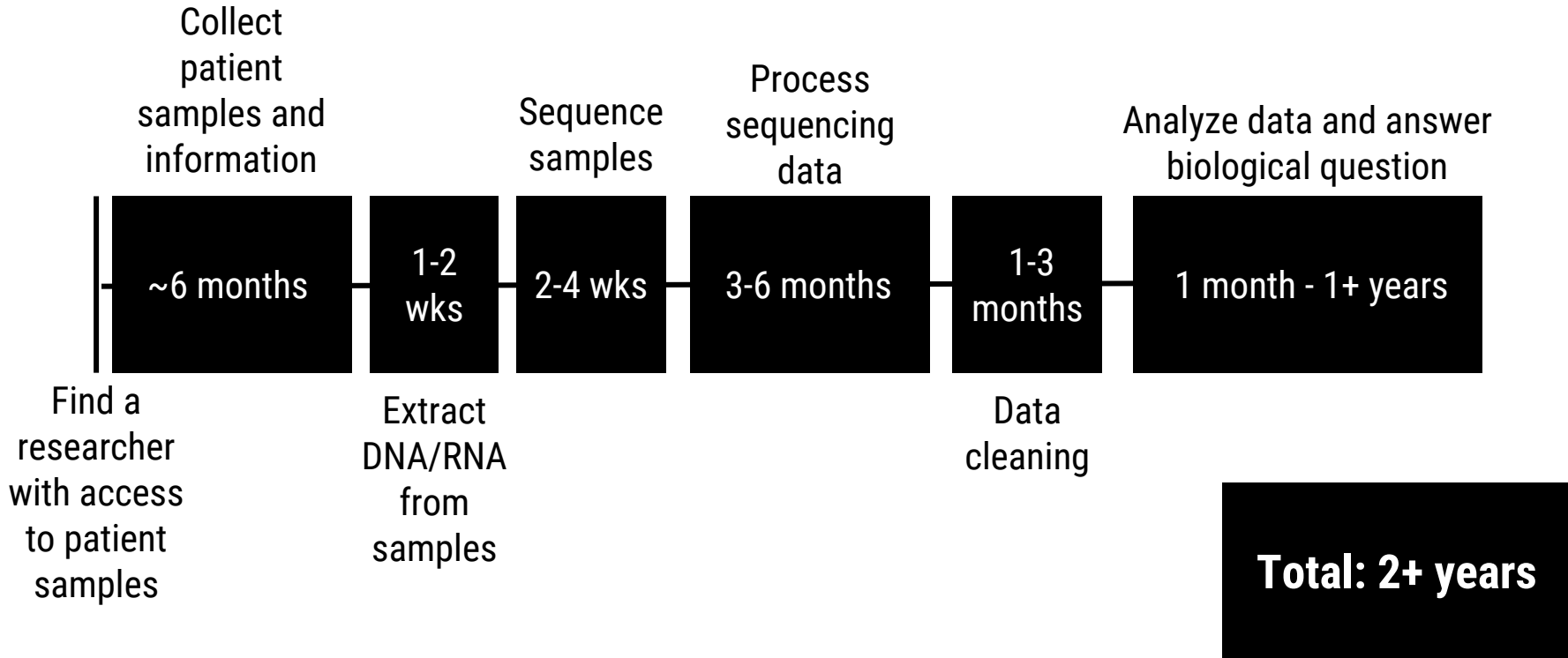


# What makes primary cancer different than metastatic cancer?



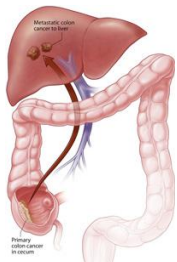


# What makes primary cancer different than metastatic cancer?

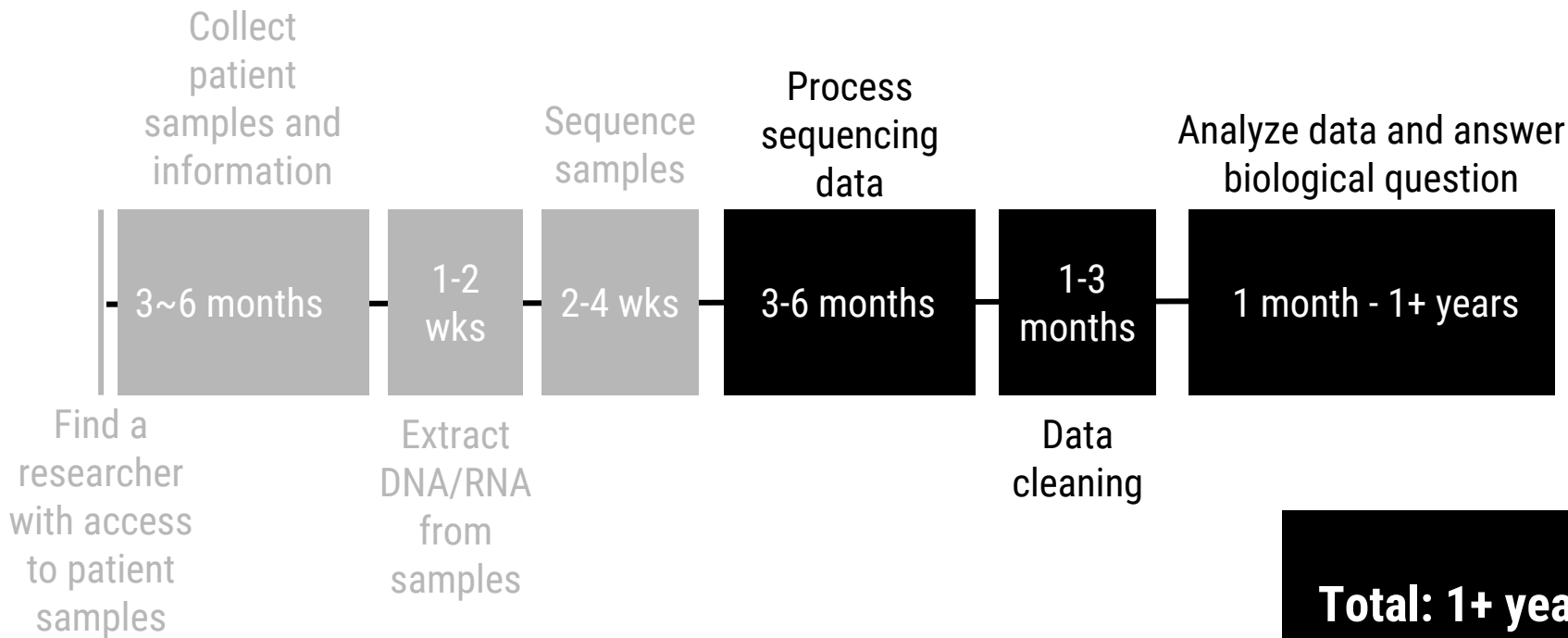




Biologists have recently gotten pretty good at making their data available to the public.



# What makes primary cancer different than metastatic cancer?

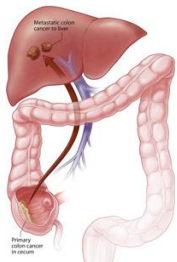




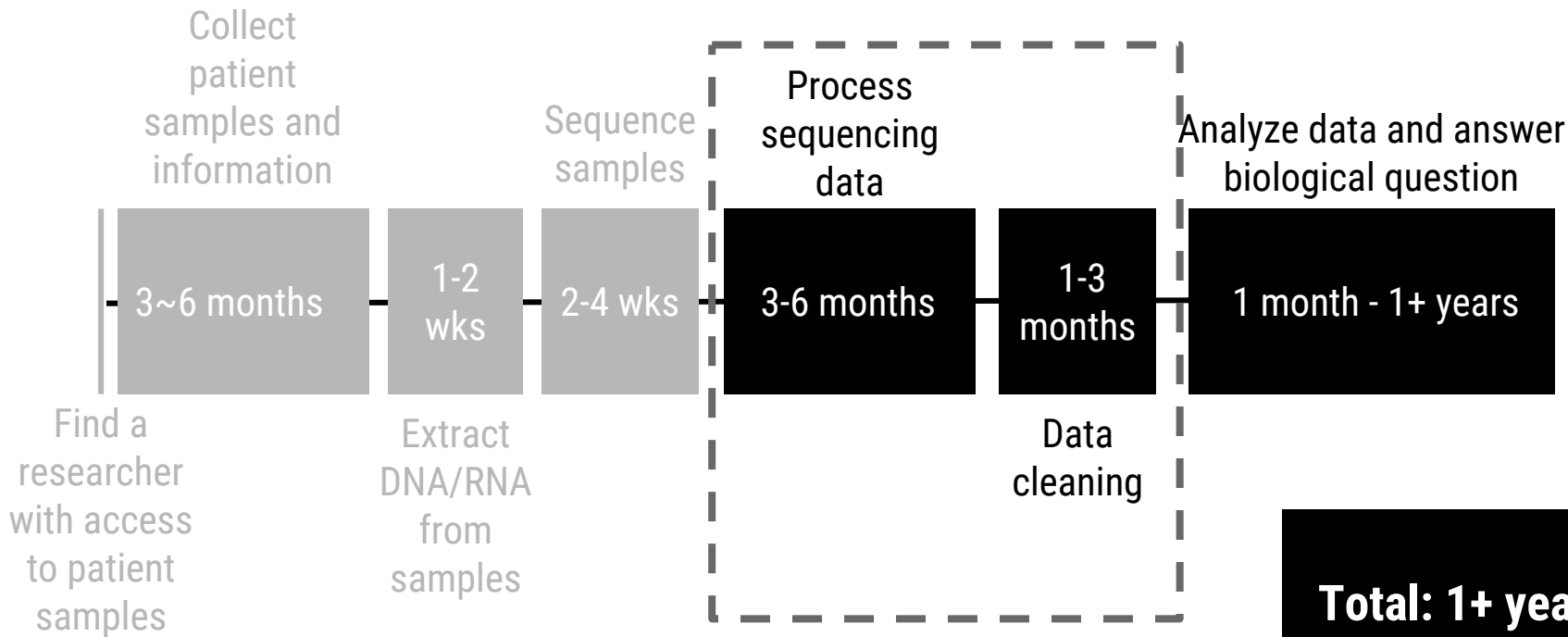
Biologists have recently gotten pretty good at making their data available to the public.



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



# What makes primary cancer different than metastatic cancer?

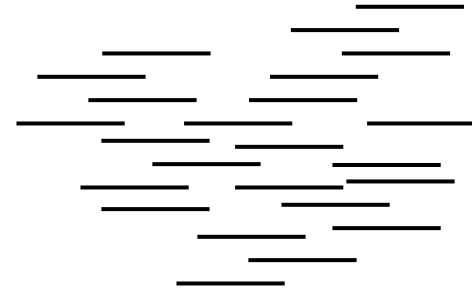
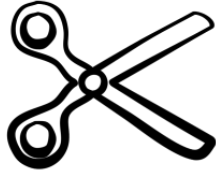
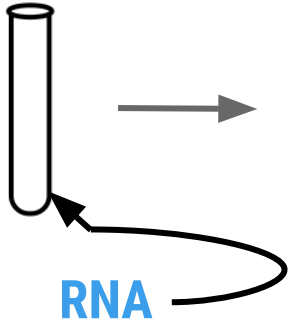




# Measuring Transcription



# Next Generation Sequencing (NGS) in one slide



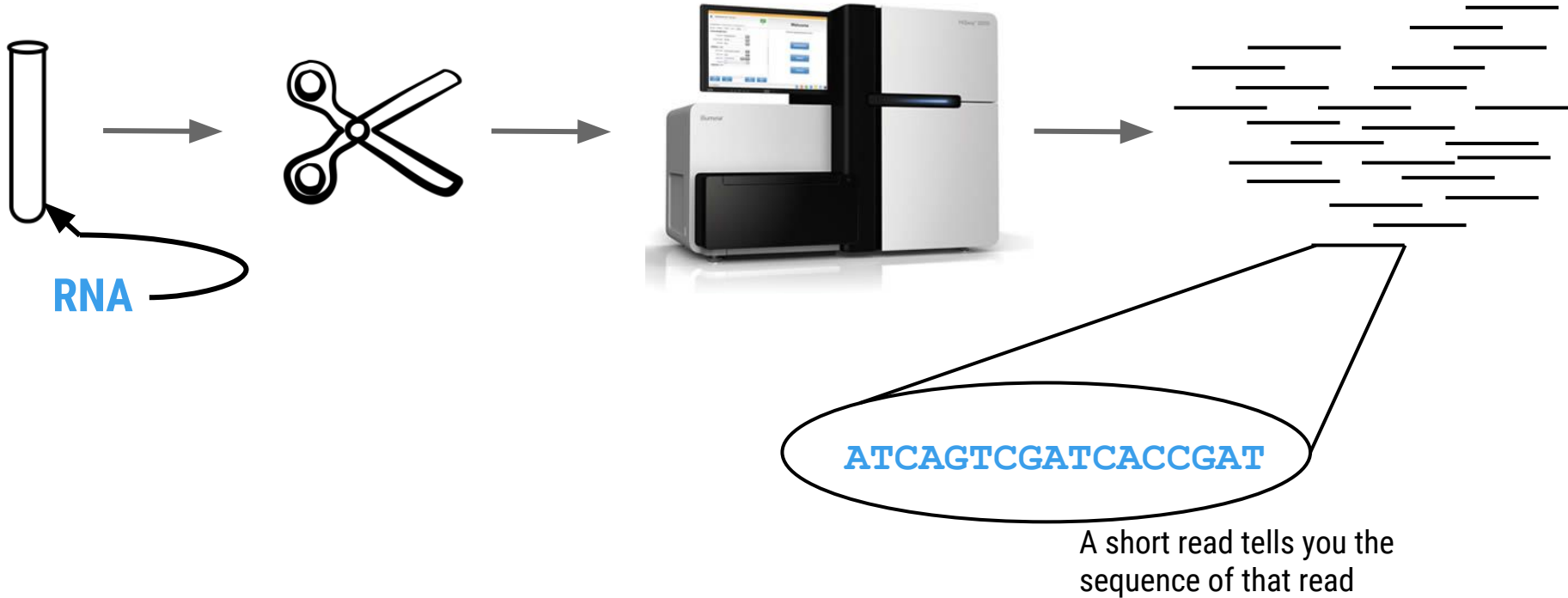
**Step 1:** Extract RNA to get sample of interest

**Step 2:** Chop up RNA into smaller pieces

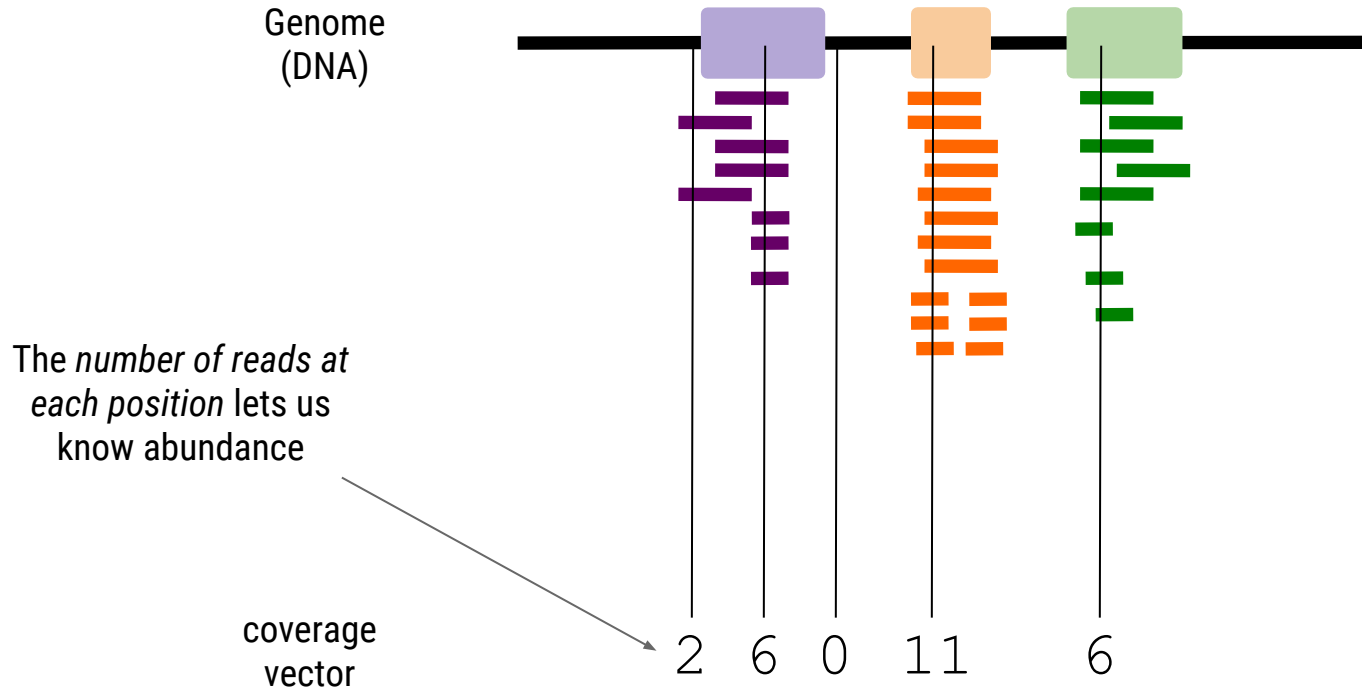
**Step 3:** Sequence the sample

**Step 3:** Obtain short read data from the sequencer

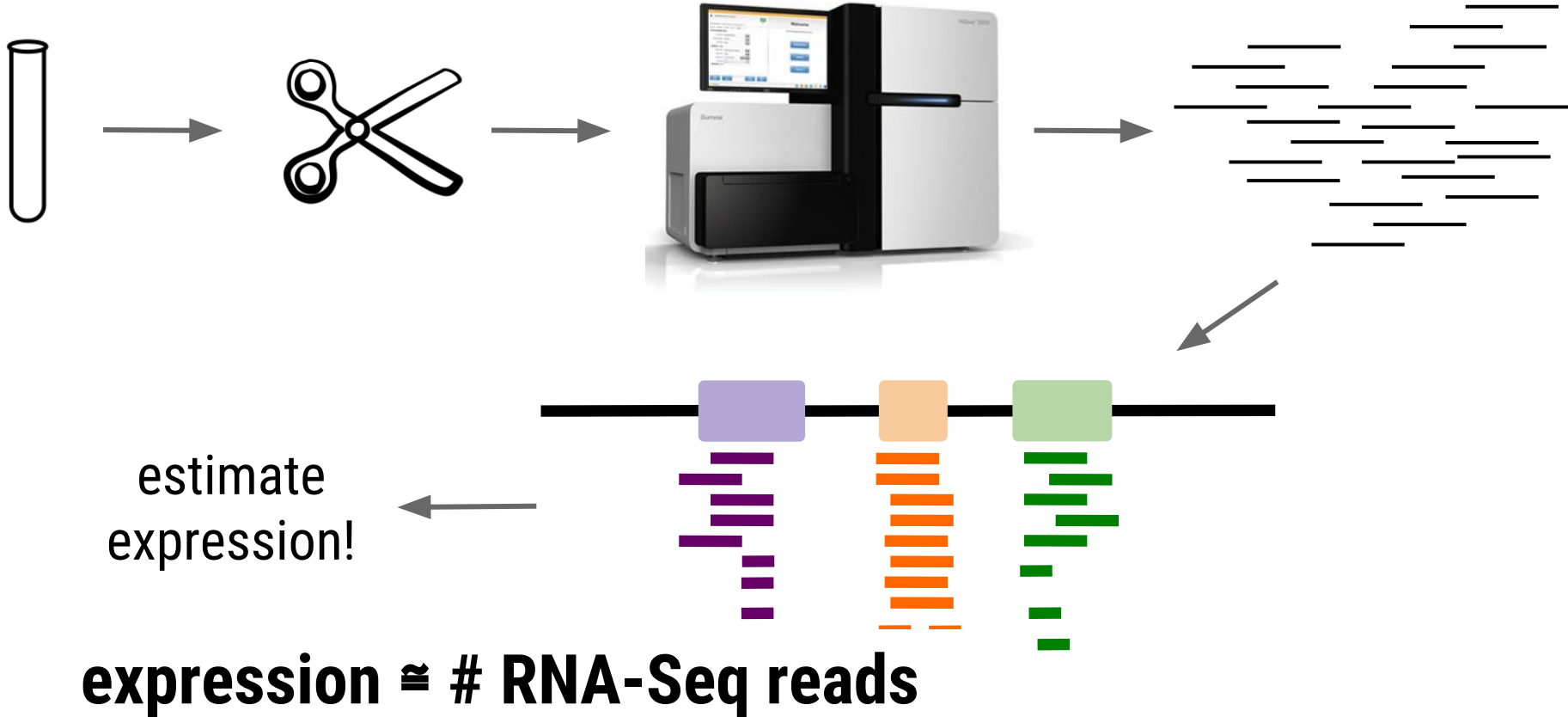
# Next Generation Sequencing (NGS) in one slide



# We first need to align these reads back to the genome

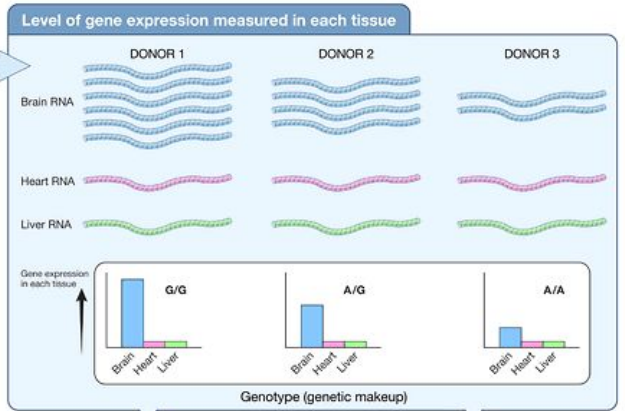
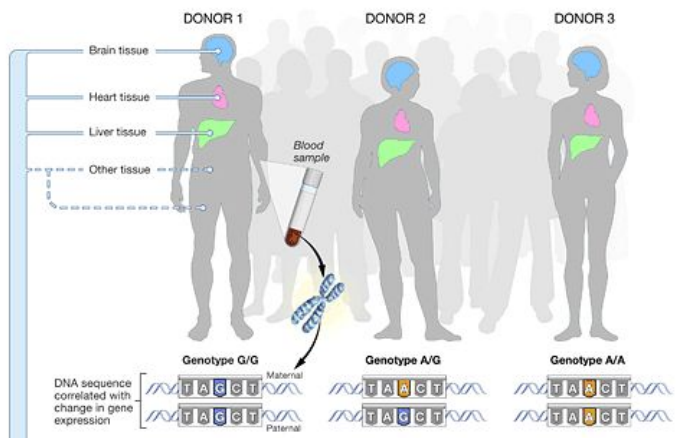


RNA-Seq = estimate expression across entire genome

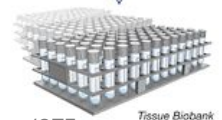


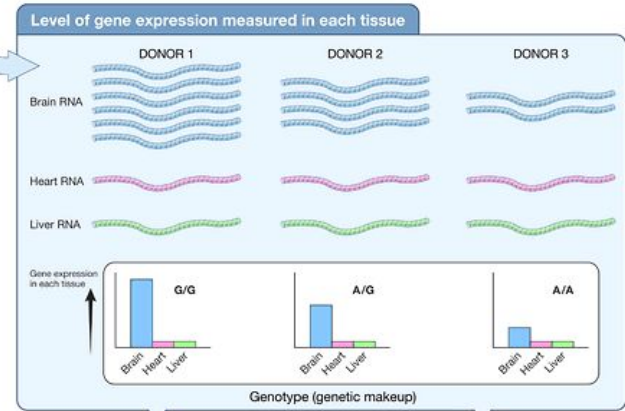
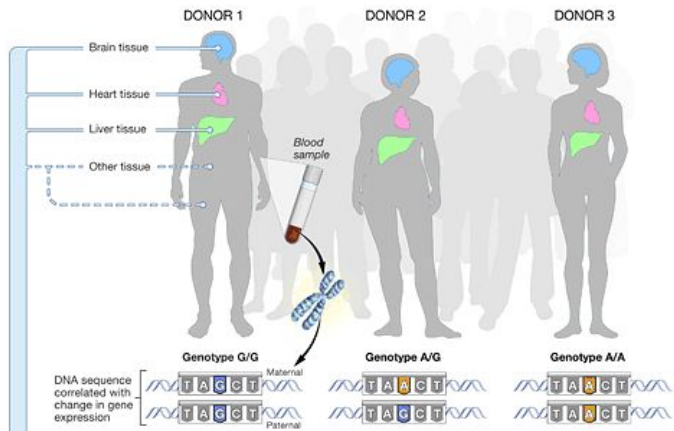
# Scaling Up



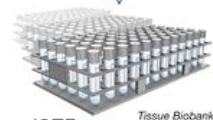


**GTEX**





**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

TCGA data describes ...including  
**33** DIFFERENT TUMOR TYPES  
**10** RARE CANCERS

To put this into perspective, **1 petabyte** of data is equal to  
**212,000** DVDs

...based on paired tumor and normal tissue sets collected from  
**11,000** PATIENTS  
...using  
**7** DIFFERENT DATA TYPES

**TCGA RESULTS & FINDINGS**

	<b>MOLECULAR BASIS OF CANCER</b>	Improved our understanding of the genomic underpinnings of cancer	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.
	<b>TUMOR SUBTYPES</b>	Revolutionized how cancer is classified	TCGA revolutionized how cancer is classified by identifying tumor subtypes with distinct sets of genomic alterations.*
	<b>THERAPEUTIC TARGETS</b>	Identified genomic characteristics of tumors that can be targeted with currently available therapies or used to help with drug development	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 types, including a new subtype characterized by infection with Epstein-Barr virus.

[www.cancer.gov/ccg](http://www.cancer.gov/ccg)



SRA

SRA

[Advanced](#) [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

<b>Project</b>	<b>No. of Sample</b>
<b>GTE<sub>x</sub></b> Genotype-Tissue Expression Project	9,962
<b>TCGA</b> The Cancer Genome Atlas	11,284
<b>SRA</b> Sequence Read Archive	49,848

# recount2

## 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, RangedSummarizeExperiment 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  entries

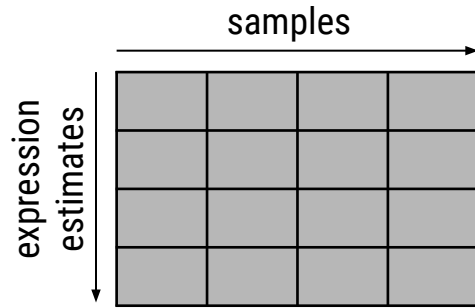
Search:

accession	number of samples	species	abstract	gene	exon	junctions	phenotype	files info
<input type="text" value="All"/>	<input type="text" value="All"/>	<input checked="" type="text" value="All"/>	<input type="text" value="All"/>	<input type="text"/>	<input type="text"/>	<input type="text" value="All"/>	<input type="text" value="All"/>	<input type="text"/>
<a href="#">SRP025982</a>	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	<a href="#">RSE counts</a>	<a href="#">RSE counts</a>	<a href="#">RSE jx_bed jx_cov counts</a>	<a href="#">link</a>	<a href="#">link</a>

The logo for 'recount2' features three vertical bars of different heights and colors: a red bar on the left, a green bar in the middle, and a blue bar on the right. The text 'recount2' is positioned to the right of these bars.

# recount2

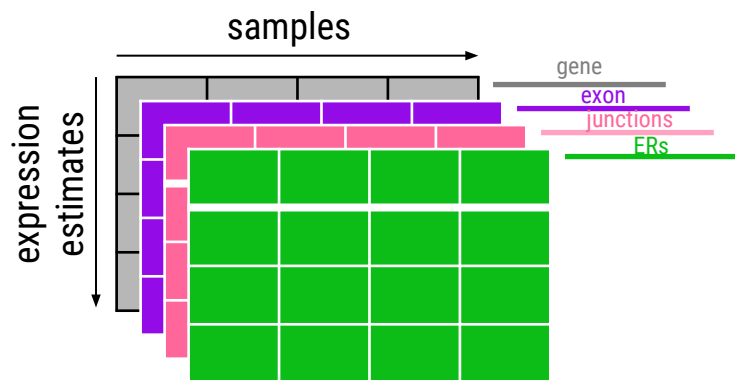
*expression data for ~70,000 human samples*



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

# recount2

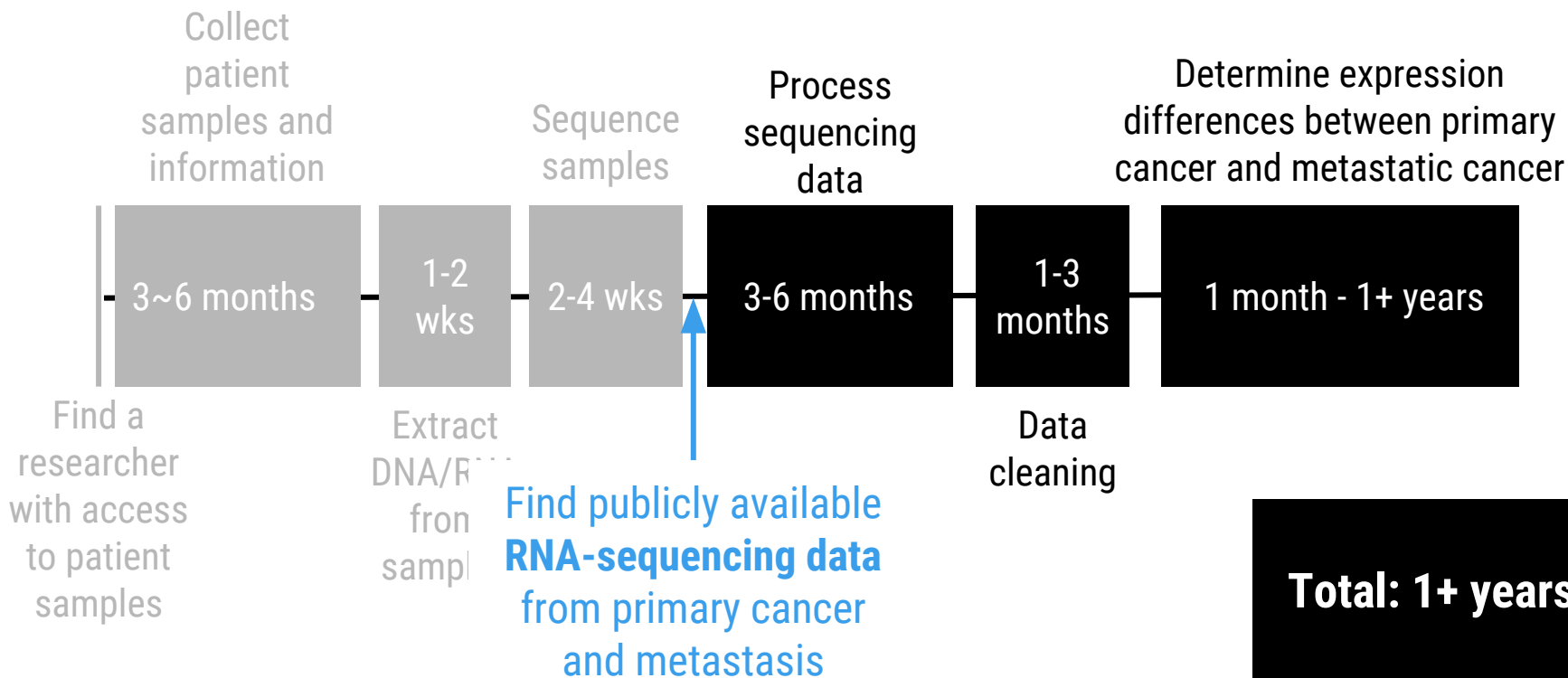
*expression data for ~70,000 human samples*

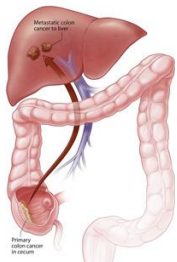


GTEEx	SRA	TCGA
N=9,962	N=49,848	N=11,284

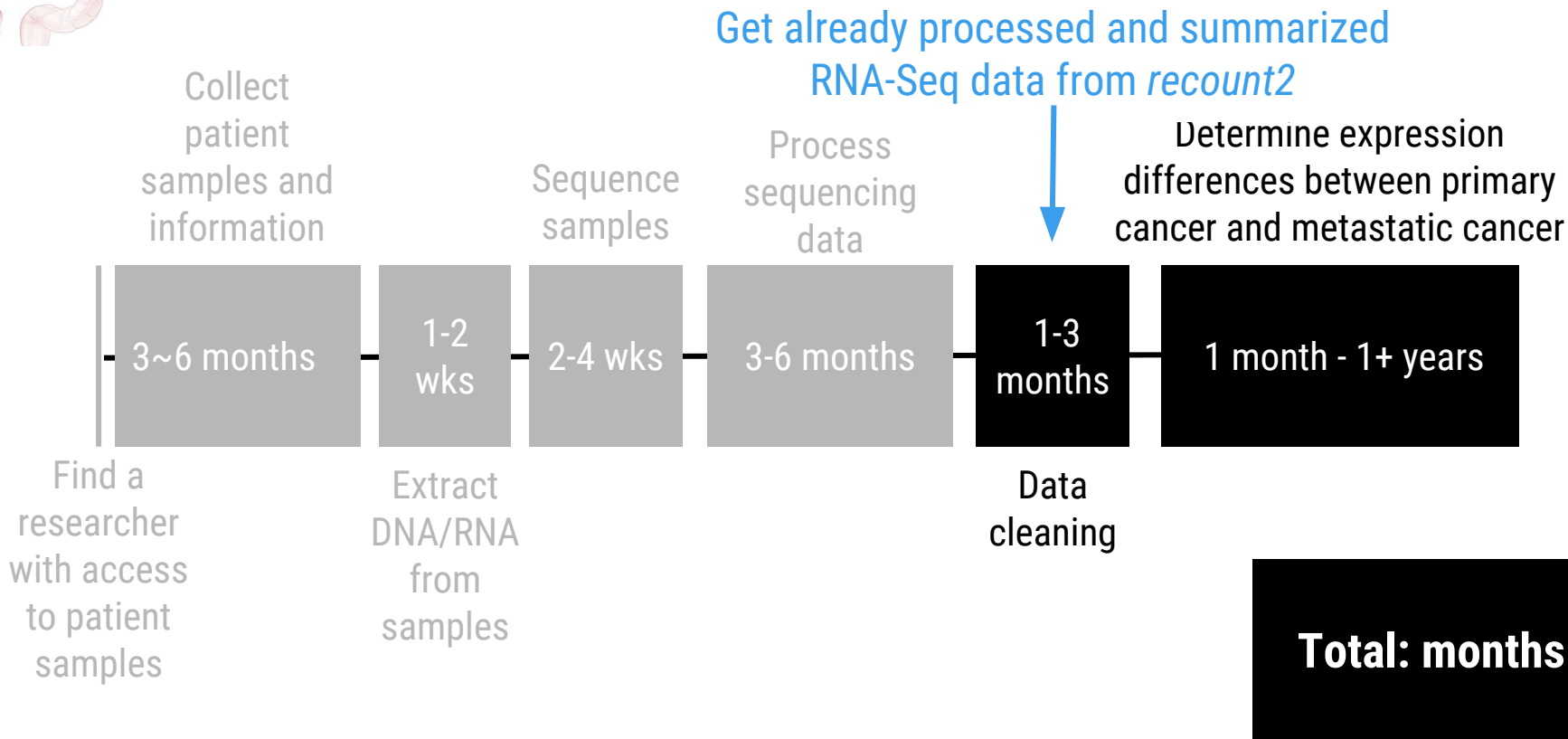


# What makes primary cancer different than metastatic cancer?



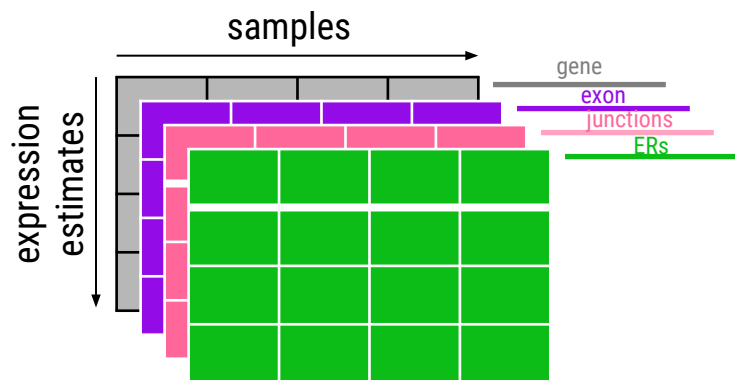


# What makes primary cancer different than metastatic cancer?



# recount2

*expression data for ~70,000 human samples*



GTEx N=9,962	SRA N=49,848	TCGA N=11,284
-----------------	-----------------	------------------

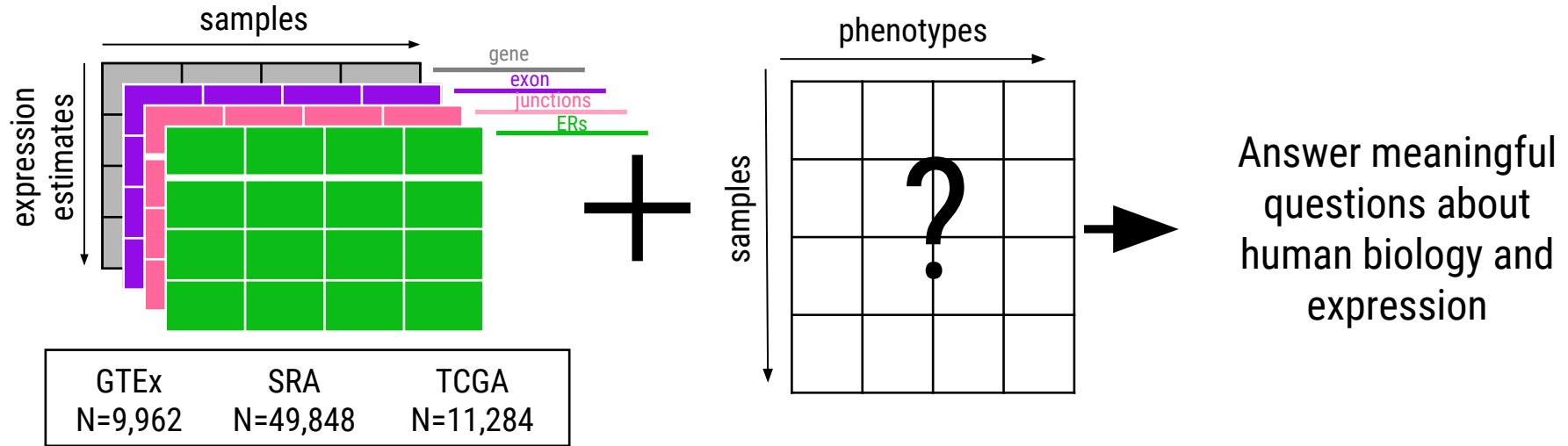


Answer meaningful  
questions about  
human biology and  
expression



# recount2

*expression data for ~70,000 human samples*



*in-silico*  
Phenotyping



# SRA phenotype information is far from complete

	<b>Sex</b>	<b>Tissue</b>	<b>Race</b>	<b>Age</b>
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 phenotype information is far from complete

	<b>Sex</b>	<b>Tissue</b>	<b>Race</b>	<b>Age</b>
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

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

Sex across the **SRA**:

<b>Level</b>	<b>Frequency</b>
F	95
female	2036
Female	51
M	77
male	1240
Male	141
<b>Total</b>	<b>3640</b>

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

**Sex across the SRA:**

<b>Level</b>	<b>Frequency</b>
F	95
female	2036
Female	51
M	77
male	1240
Male	141
<b>Total</b>	<b>3640</b>

"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...

**Sex across the SRA:**

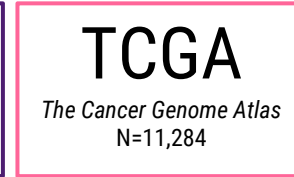
<b>Level</b>	<b>Frequency</b>
F	95
female	2036
Female	51
M	77
male	1240
Male	141
<b>Total</b>	<b>3640</b>

"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"

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

**Goal :**

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





**Goal :**

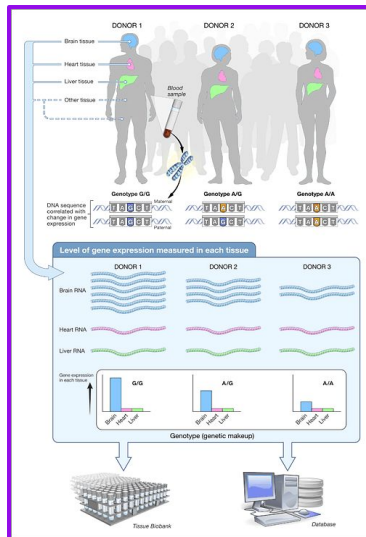
to accurately predict critical phenotype information for all samples in *recount2*



**GTE<sub>x</sub>**  
*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 analyzed over **2.5 PETABYTES** of data.

To put this into perspective, 2 petabyte of data is equal to:

**212,000 DVD's**

including **33 DIFFERENT TUMOR TYPES** and **10 RARE CANCERS**.

based on paired tumor and normal tissue sets collected from **11,000 PATIENTS** using **7 DIFFERENT DATA TYPES**.

**TCGA RESULTS & FINDINGS**

- MELEIOL TUMOR CANCER**: Improved our understanding of the genomic architecture of cancer.
- TUMOR SUBTYPES**: Re-discovered new cancer subtypes.
- BIOMARKERS**: Identified potential characterization of tumors that can be targeted with currently available therapies to lead to long-term drug development.

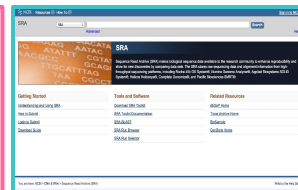
**WHAT'S NEXT?**

TCGA's identification of tumor cell genetic alterations in lung squamous cell carcinoma led to NCI's LungMap. The atlas will benefit patients based on the specific genomic changes in their tumor.

The Cancer Genome Atlas (CGA) features TCGA and other NCI genomics data sets on a common data access platform. The CGA also has many research capabilities that will allow researchers to answer more clinically-relevant questions with increased ease.

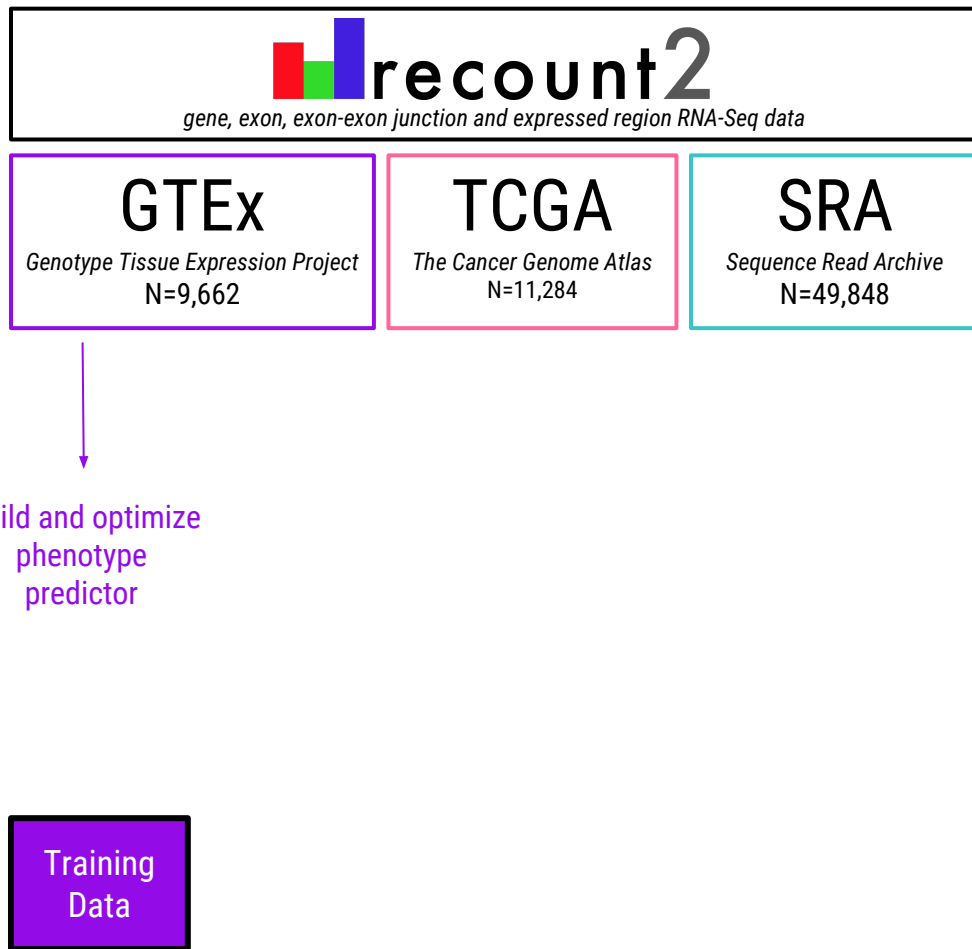
**20 COLLABORATING INSTITUTIONS** across the United States and Canada.

[www.cancer.gov/tcga](http://www.cancer.gov/tcga)



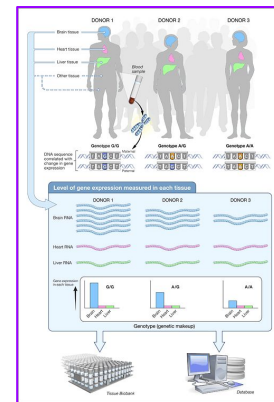
**Goal :**

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



# Missingness limited in **GTEx** phenotype data

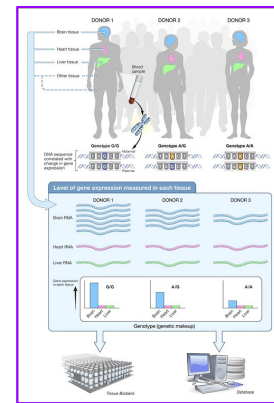
	<b>Sex</b>	<b>Tissue</b>	<b>Race</b>	<b>Age</b>
1	male	Lung	White	59
2	male	Brain	White	27
3	female	Heart	Black or African American	23
4	male	Brain	White	51
5	male	Skin	White	27
6	male	Lung	White	68
7	female	Brain	White	61
8	female	Adipose Tissue	White	42
9	male	Brain	White	40
10	female	Uterus	White	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

	<b>Sex</b>	<b>Tissue</b>	<b>Race</b>	<b>Age</b>
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



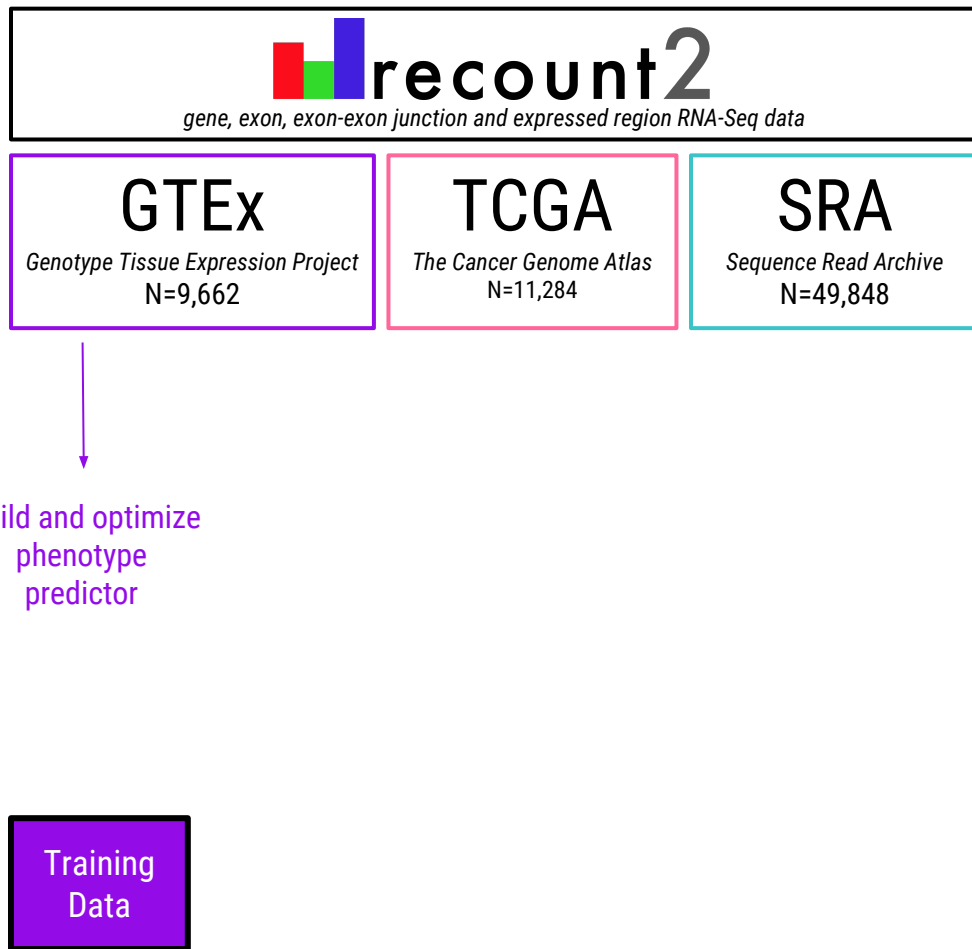
**GTEx**

## **Sex across GTEx:**

<b>level</b>	<b>Frequency</b>
female	3,626
male	6,036
NA	0

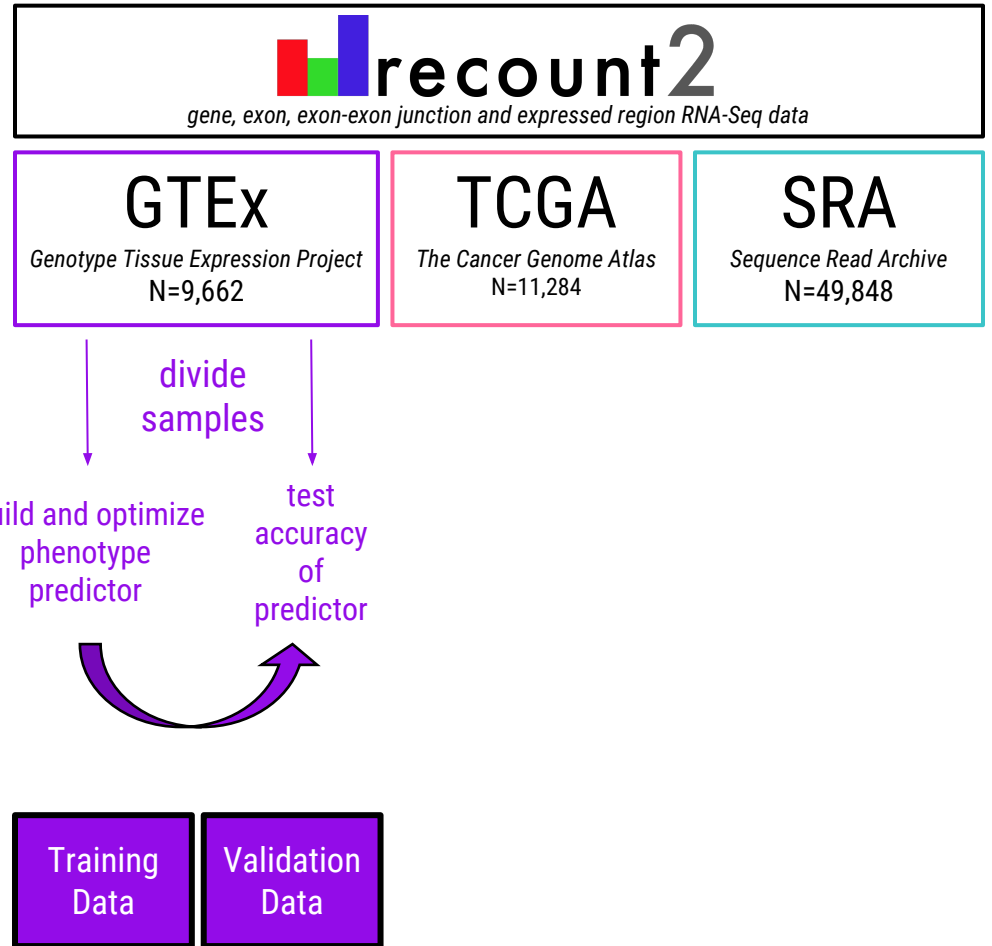
**Goal :**

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

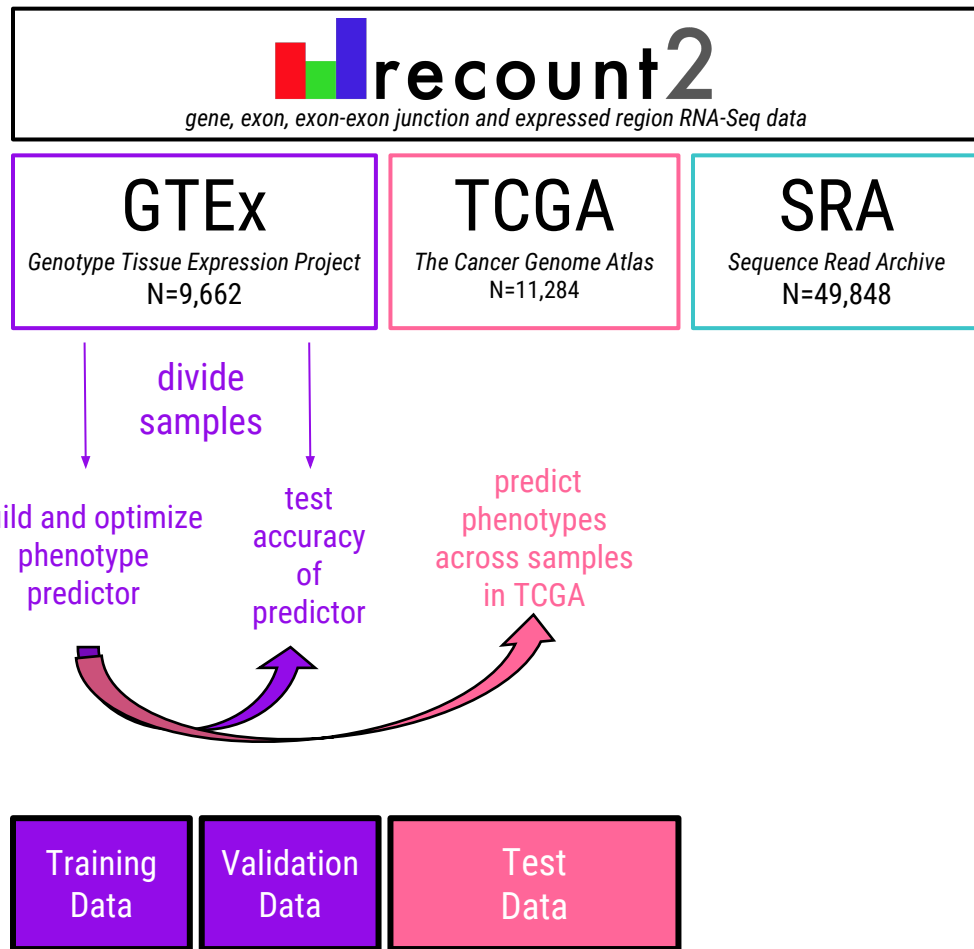


**Goal :**

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

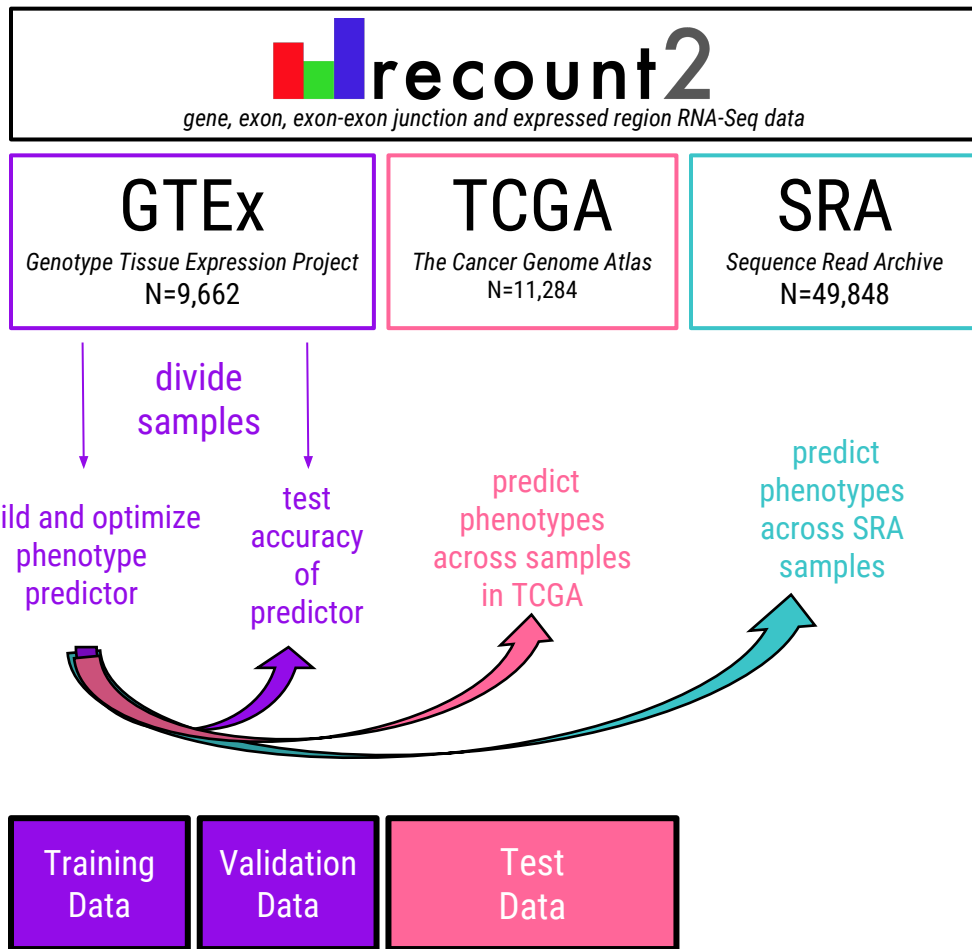


**Goal :**  
to accurately  
predict critical  
phenotype  
information for  
all samples in  
*recount2*



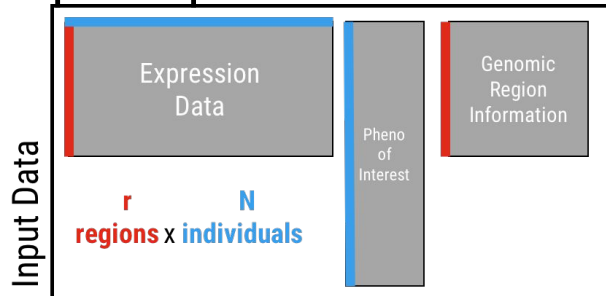
**Goal :**

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





# phenopredict



functions

- filter\_regions()
- build\_predictor()
- test\_predictor()
- extract\_data()
- predict\_pheno()

test\_predictor()

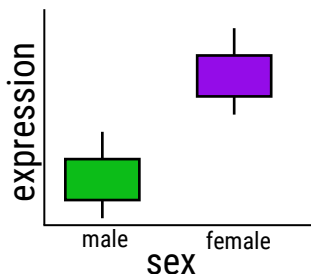
Predict phenotype and assess accuracy in training set data

predictions	reported
male	male
female	female
female	female
male	male

Prediction accuracy: 100%

filter\_regions()

Identify regions with differential expression for each level



extract\_data()

Extract expression information at regions identified by filter\_regions() in a new data set

new data set samples

expression @ filtered regions				

build\_predictor()

Extract coefficient estimates across regions

expression ~ phenotype

$\beta_r$  phenotype ( $P_i$ )

filtered regions (r)	male	female

predict\_pheno()

Predict phenotypes across samples in this new data set

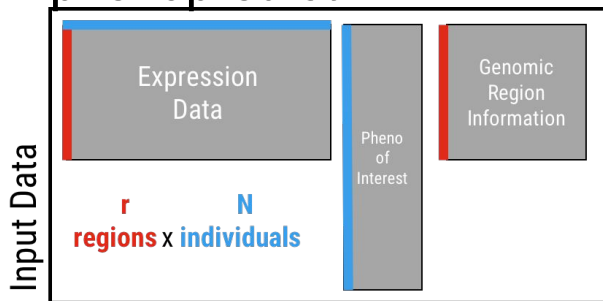
apply coefficient estimates to the extracted data

predictions in new data set

male
male
male
female

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

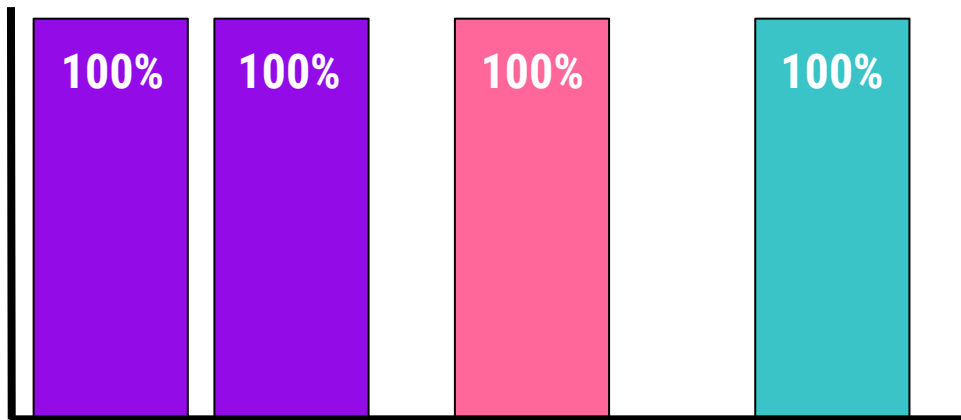
# phenopredict



functions

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

Accuracy



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

Training  
Data

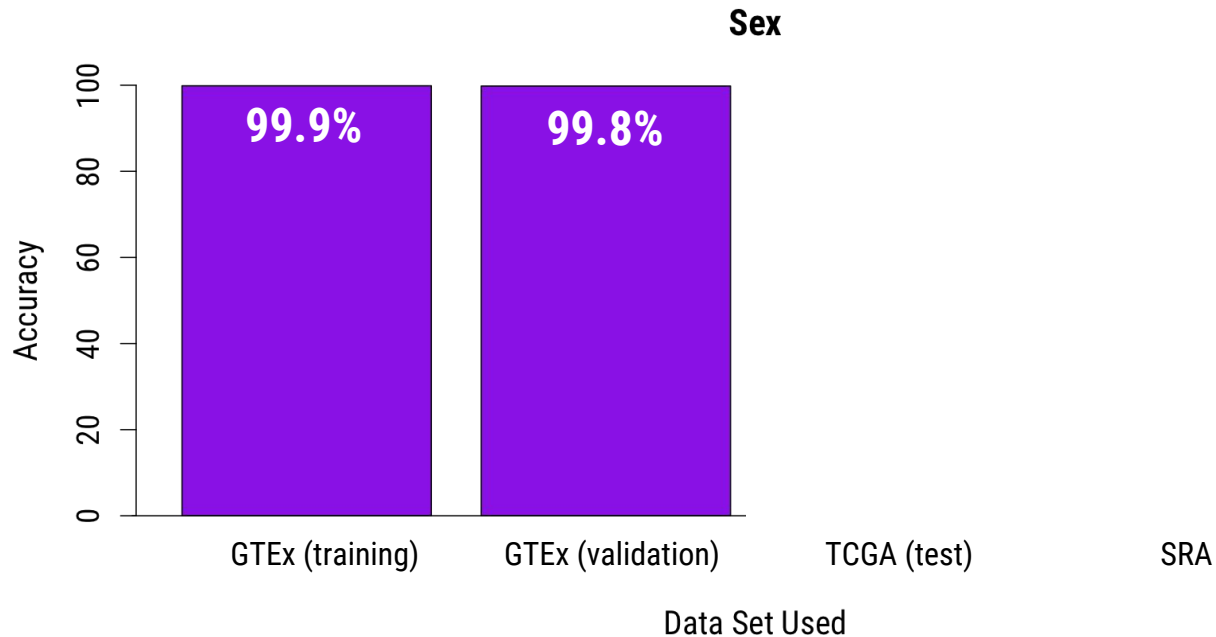
Validation  
Data

Test  
Data

Make  
predictions!

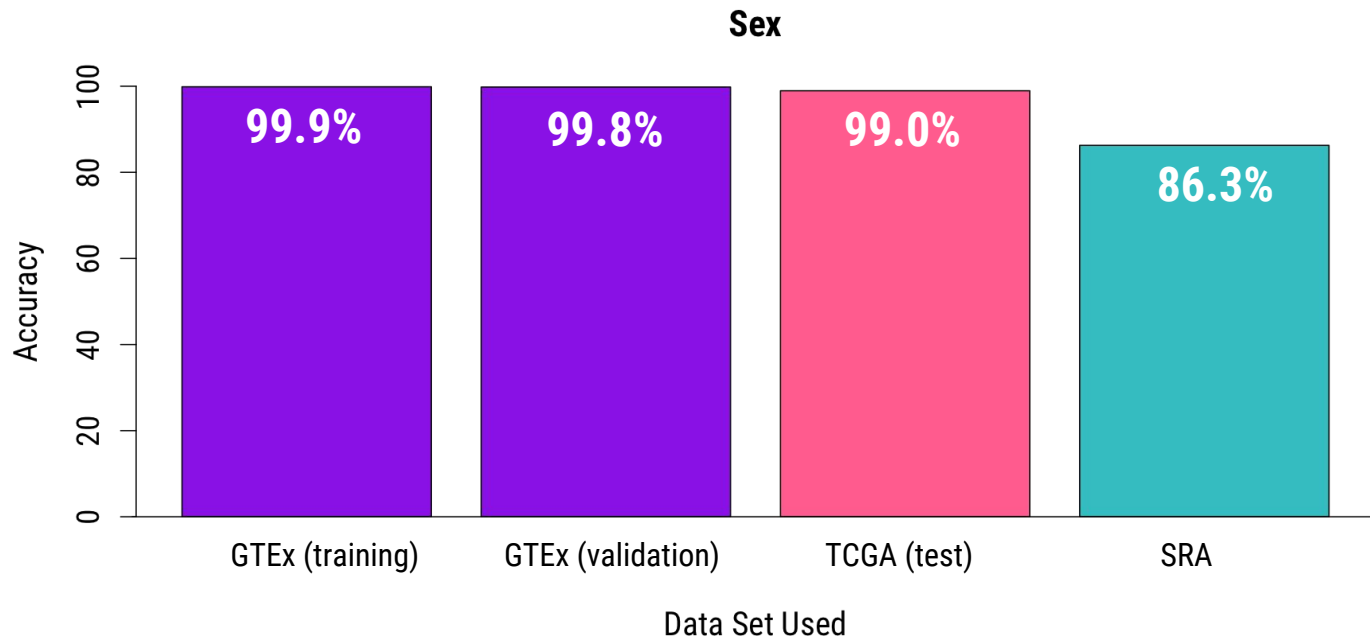
Let's get predicting...

**Sex  
prediction is  
accurate  
across data  
sets**



<b>Number of Regions</b>	<b>40</b>	<b>40</b>
<b>Number of Samples (N)</b>	<b>4,769</b>	<b>4,769</b>

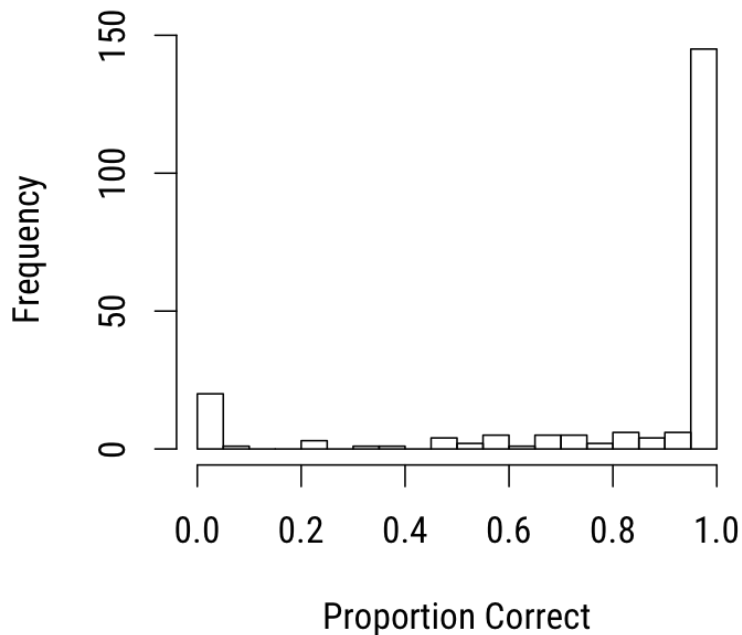
**Sex  
prediction is  
accurate  
across data  
sets**



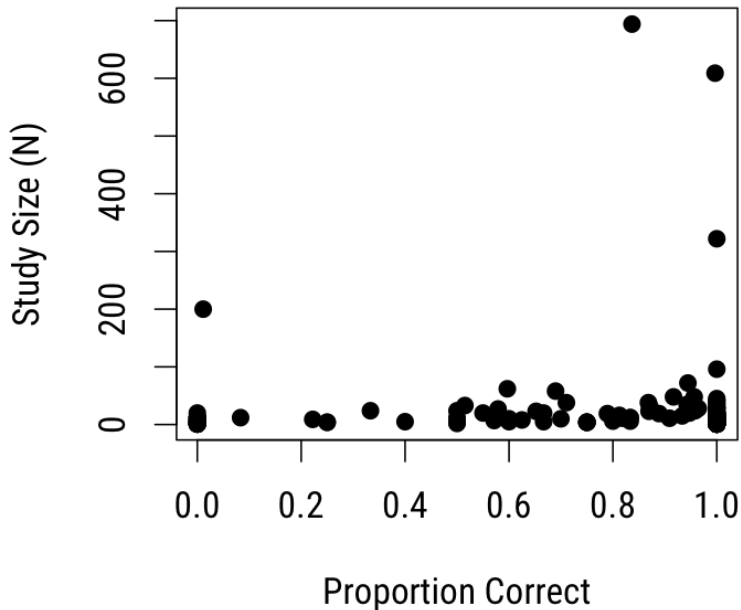
<b>Number of Regions</b>	<b>40</b>	<b>40</b>	<b>40</b>	<b>40</b>
<b>Number of Samples (N)</b>	<b>4,769</b>	<b>4,769</b>	<b>11,245</b>	<b>3,640</b>

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

**Prediction Accuracy by Study: Sex**



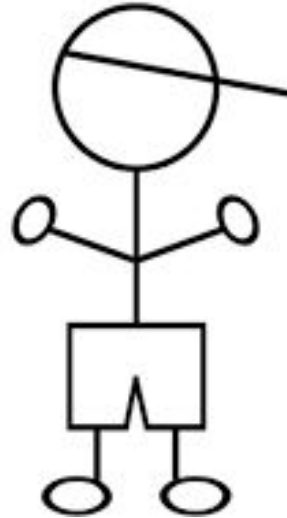
**Study Size vs. Prediction Accuracy**



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



XX

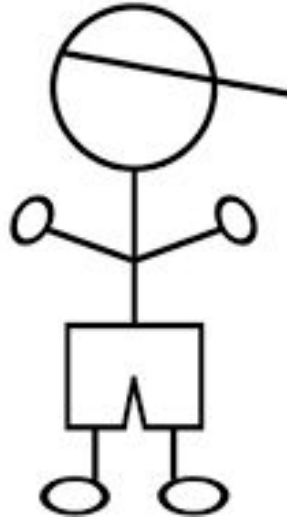


XY

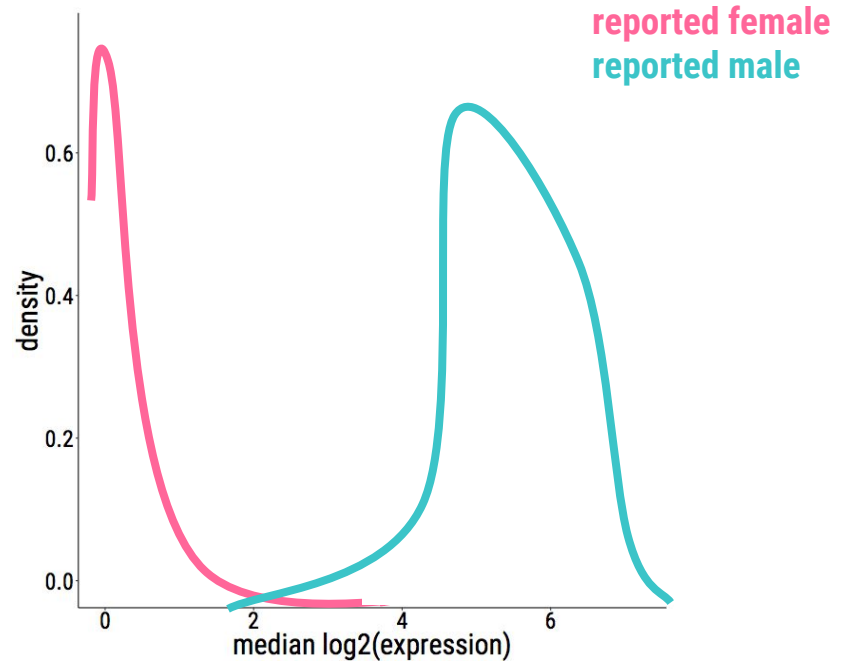
To assess misreporting of sex in the SRA, we can use Y-chromosome expression



XX

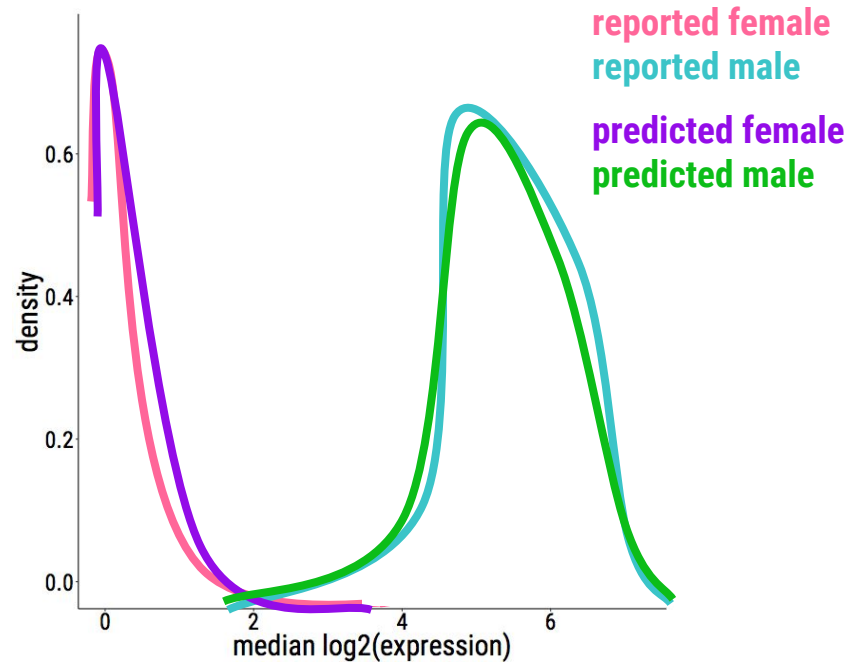
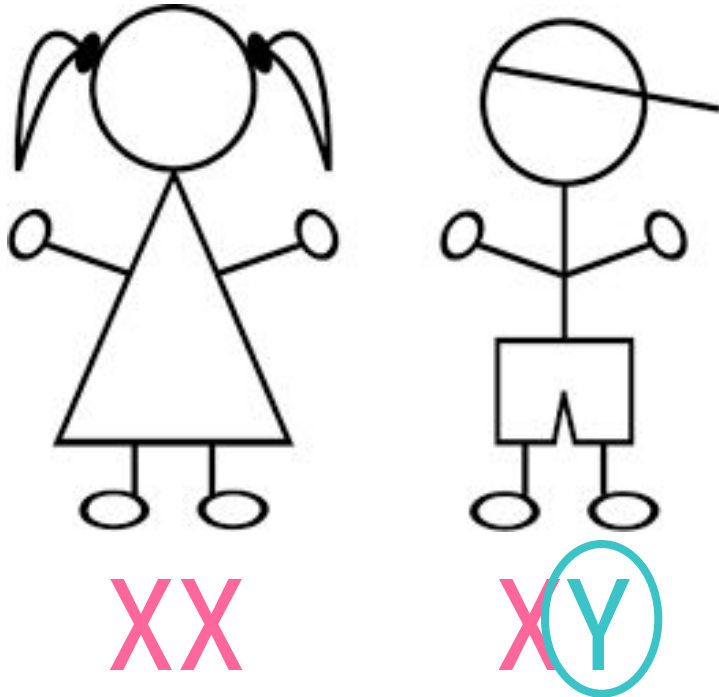


X~~Y~~



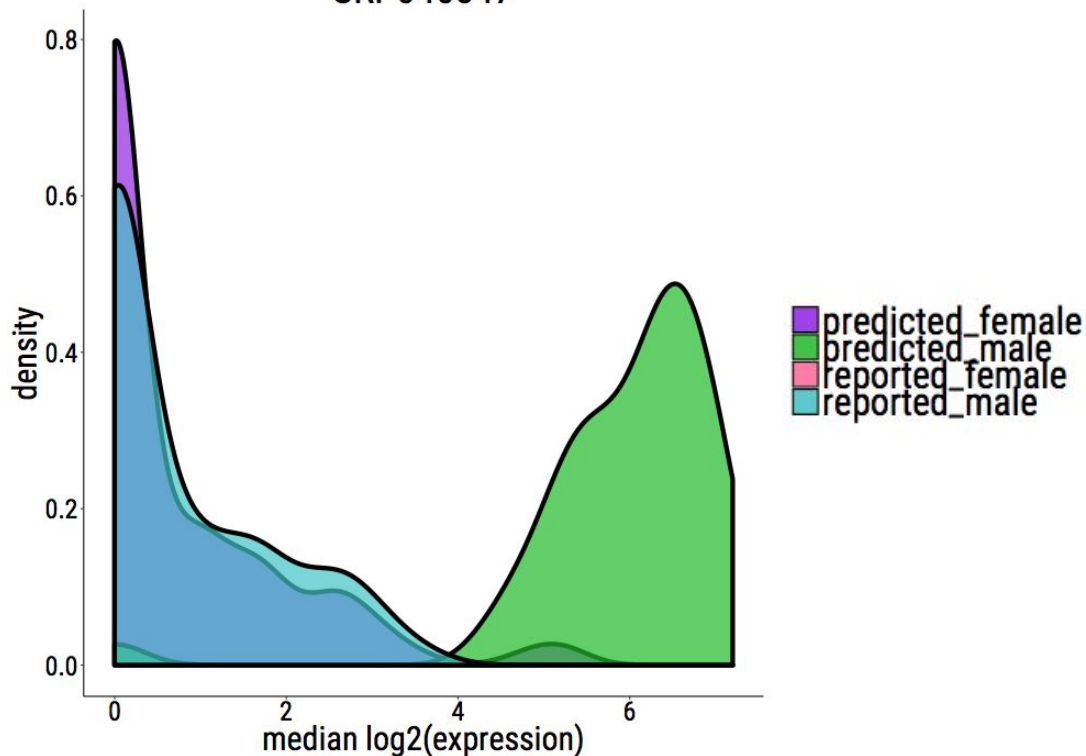


# To assess misreporting of sex in the SRA, we can use Y-chromosome expression

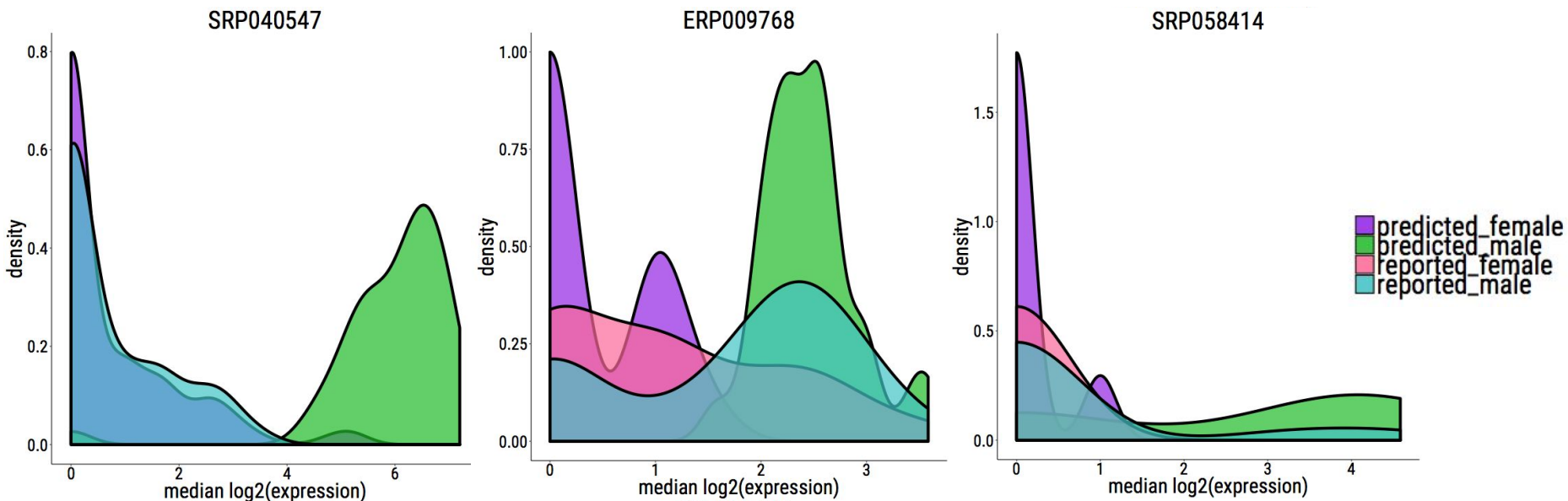


# Expression from the Y chromosome suggests misreporting of sex in the SRA

SRP040547



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

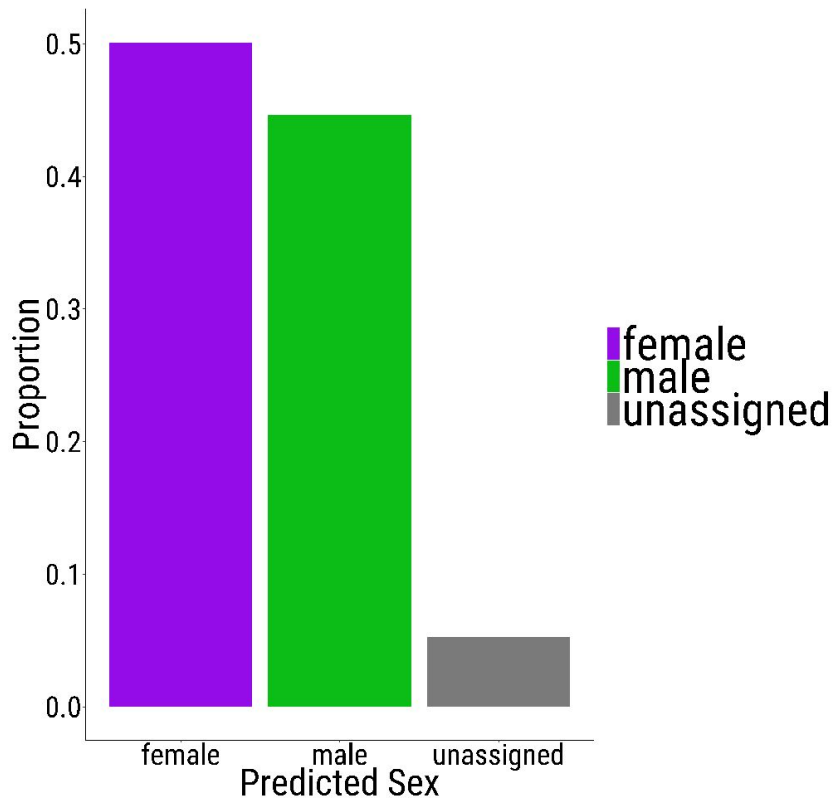


There's a well-documented history of male sex-bias in biomedical research...

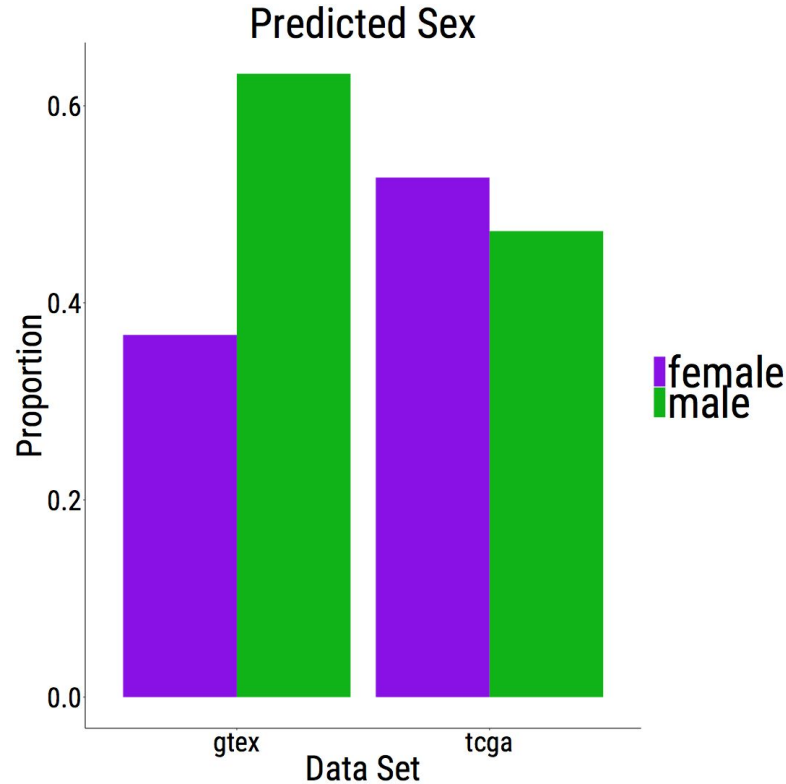
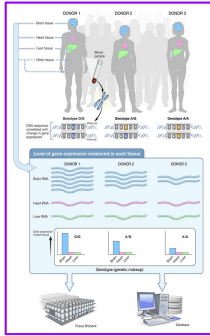


...so let's take a closer look within *recount2*

Across the ~70,000 samples in *recount2*, there are more samples predicted to be **female** than **male**.



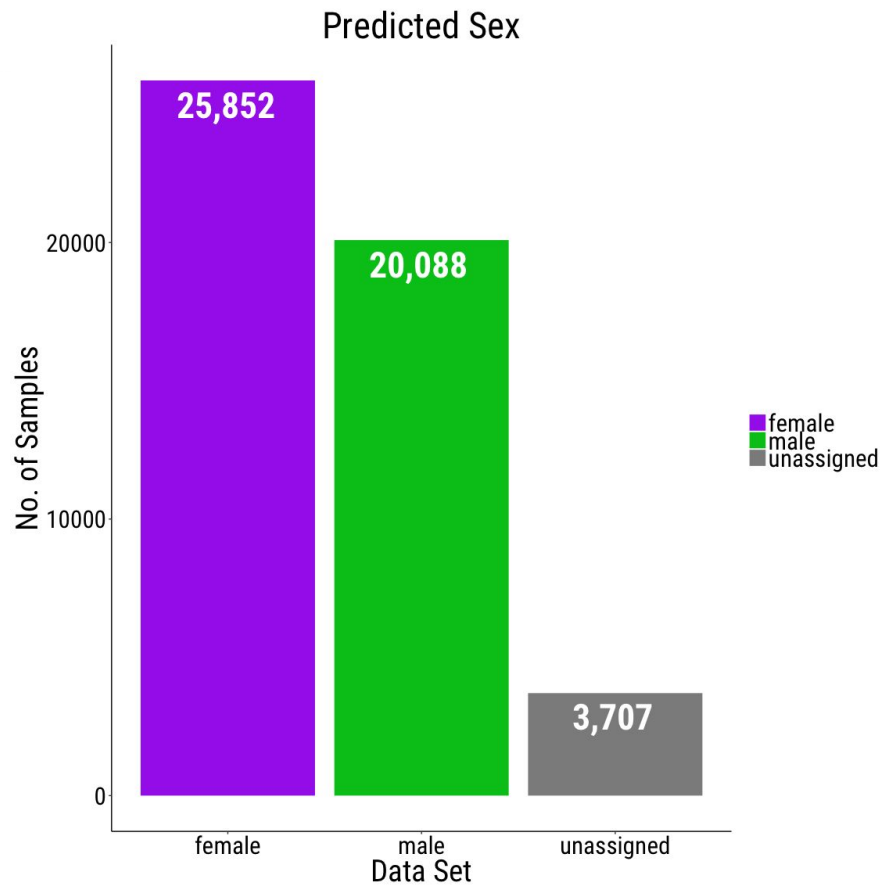
# GTEX is male-biased ; TCGA is female-biased



...but that has been previously reported by the consortia

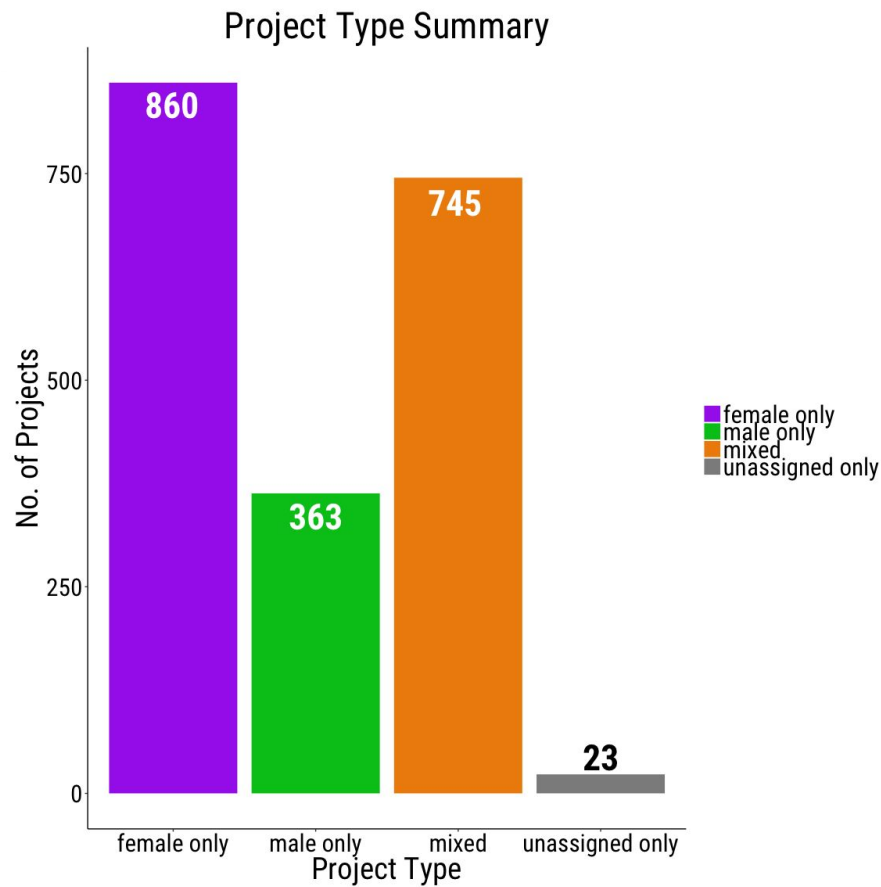
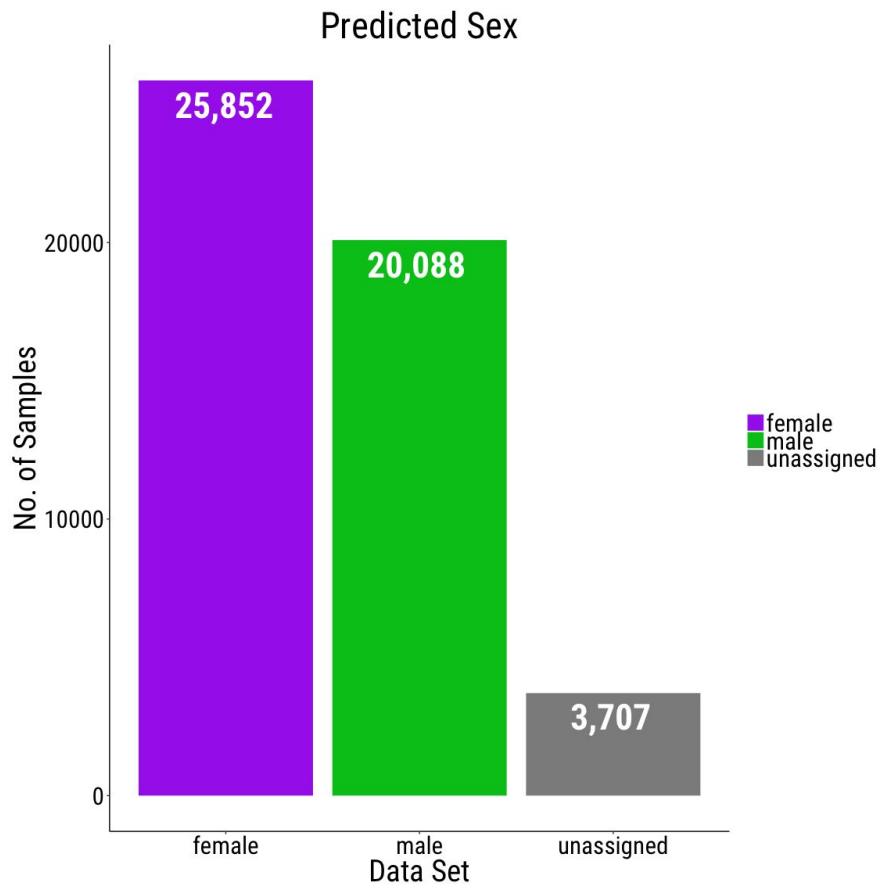
Is there a sex  
bias across  
the SRA?

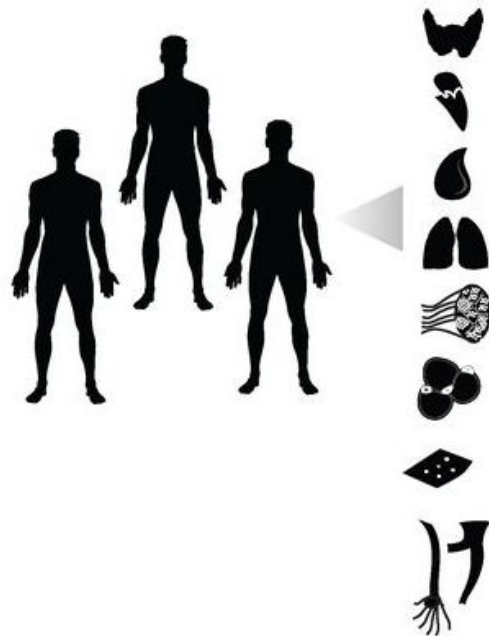
# No evidence for sex bias in samples across the SRA





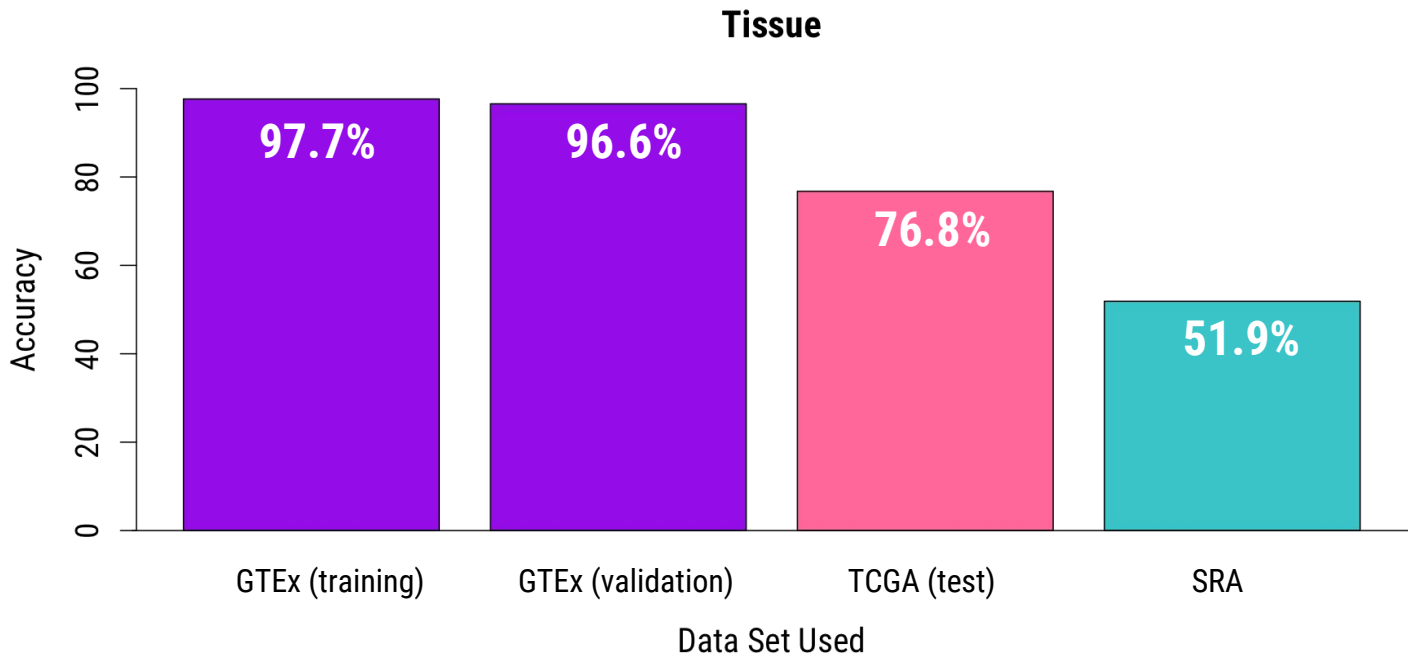
# There are many female-only, male-only, and male-female projects available in *recount*





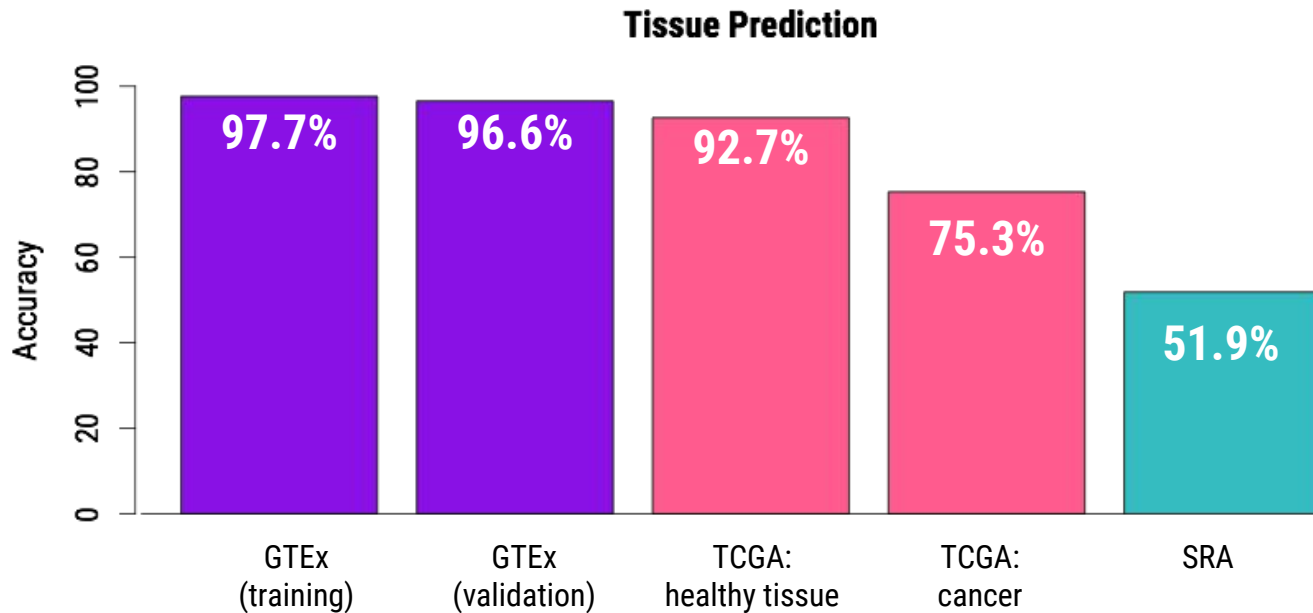
Can we use  
expression data  
to predict  
tissue?

Tissue prediction is accurate across data sets



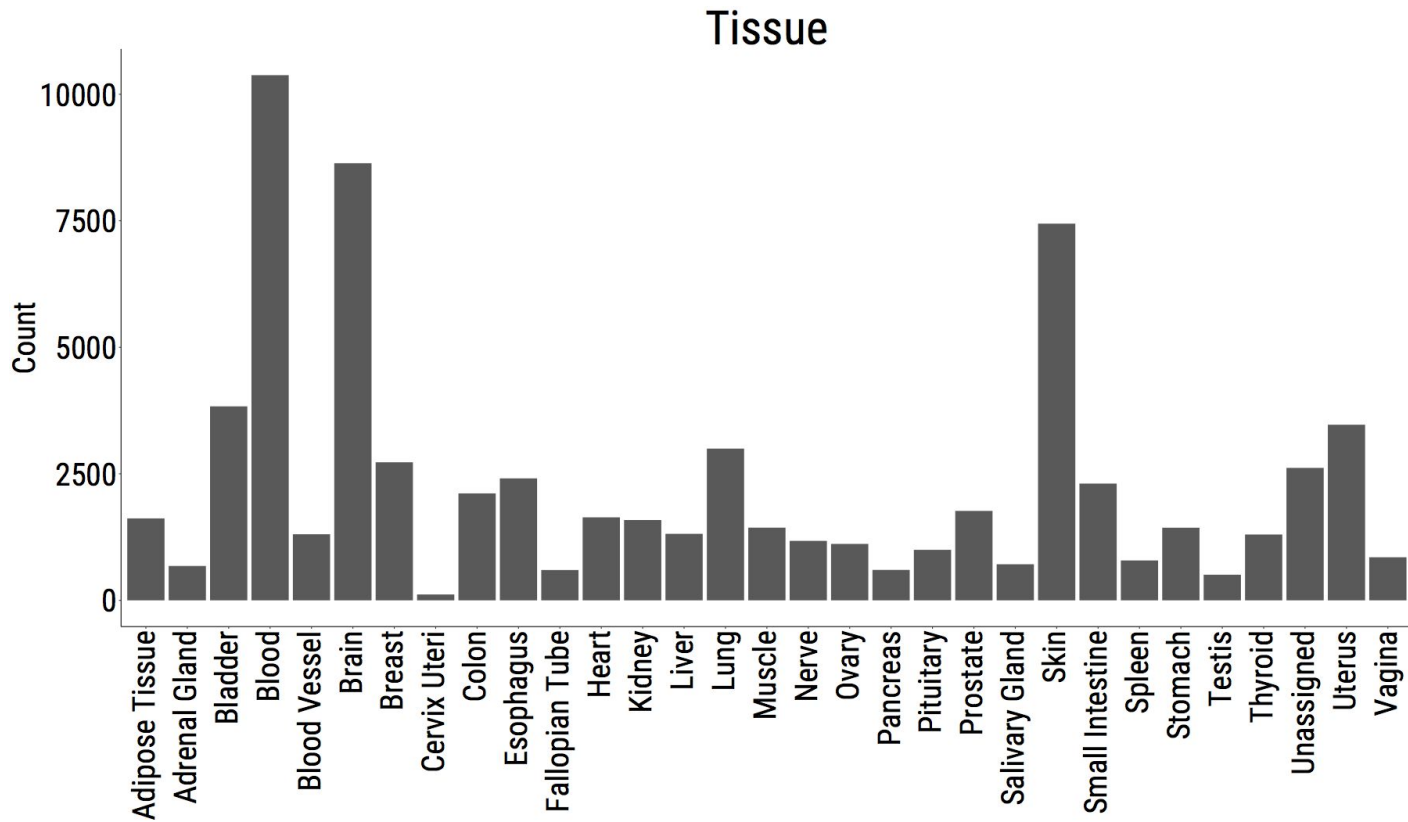
<b>Number of Regions</b>	2,281	2,281	2,281	2,281
<b>Number of Samples (N)</b>	4,769	4,769	7,317	8,951

Prediction is more accurate in healthy tissue



<b>Number of Regions</b>	2,281	2,281	2,281	2,281	2,281
<b>Number of Samples (N)</b>	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 prediction is largely accurate across *recount2*

Tissue can be accurately predicted from expression data.

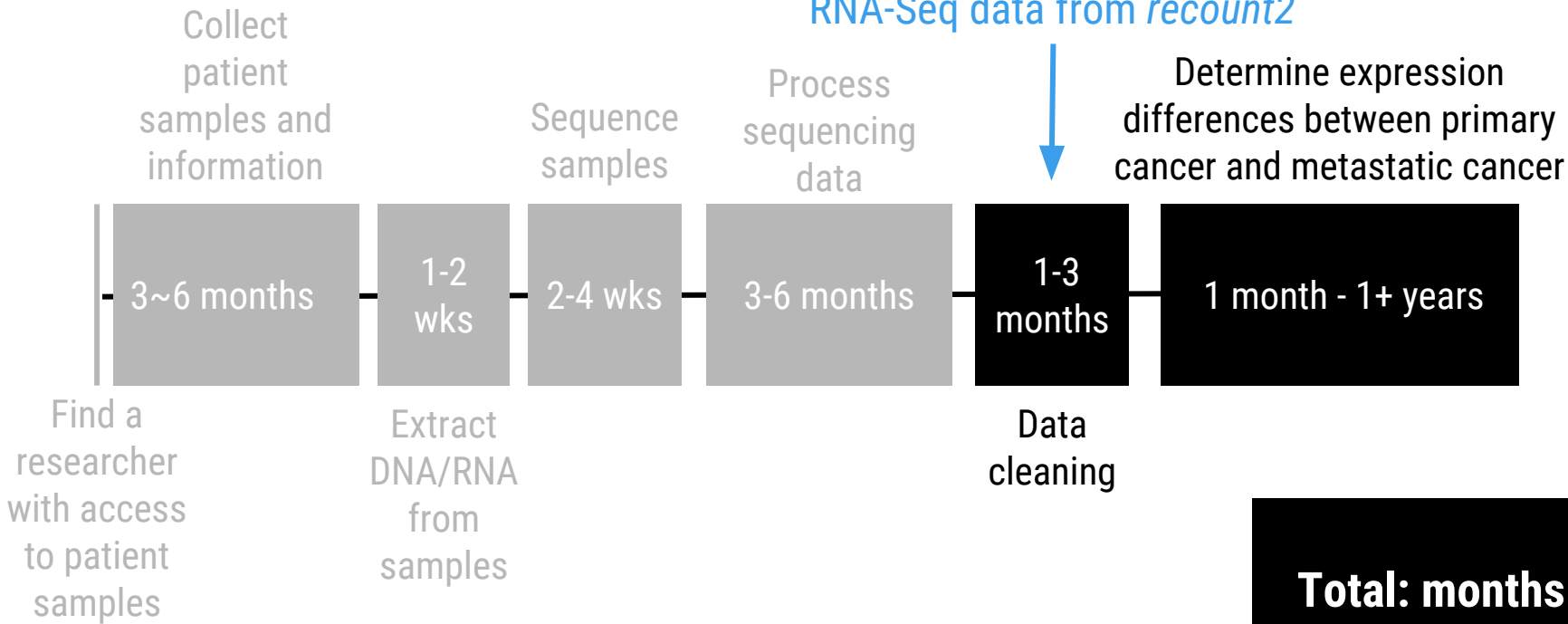
Discordant predictions are often made to biologically similar tissues.

Sometimes, predictions are inaccurate.



# What makes primary cancer different than metastatic cancer?

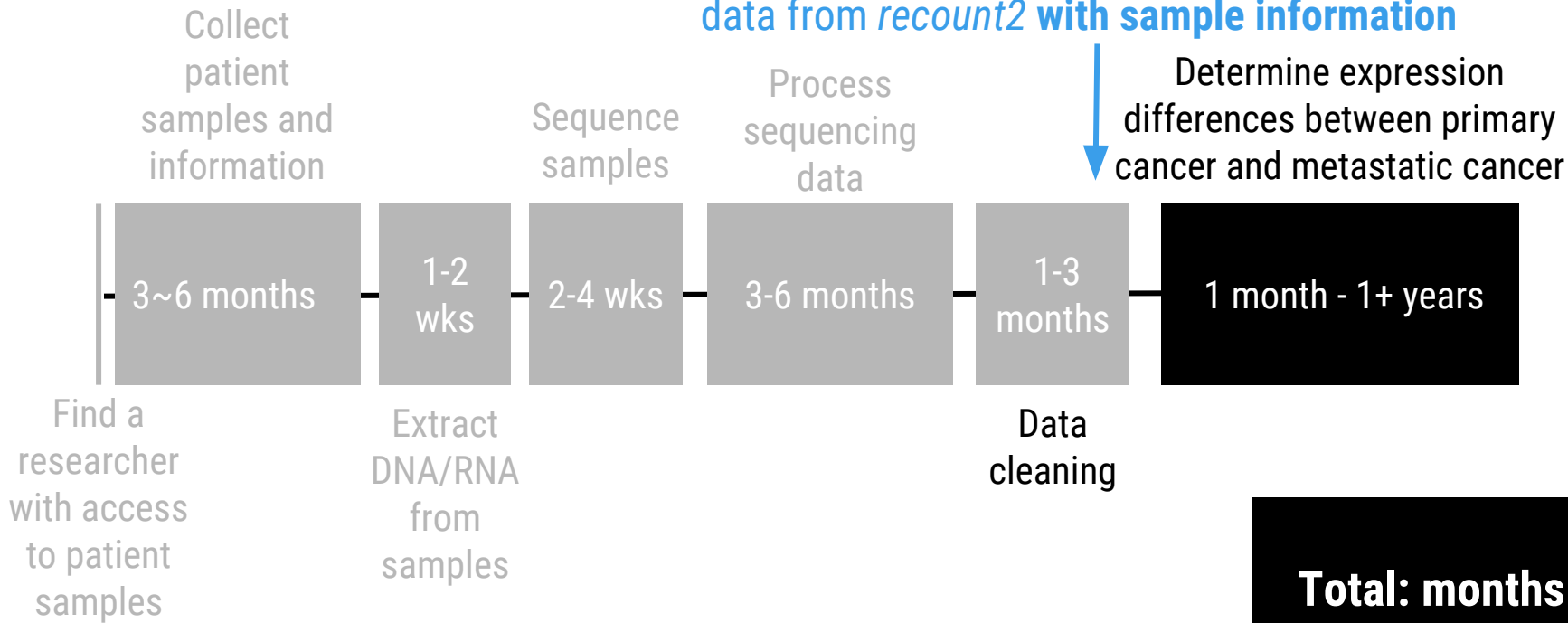
Get already processed and summarized  
RNA-Seq data from *recount2*





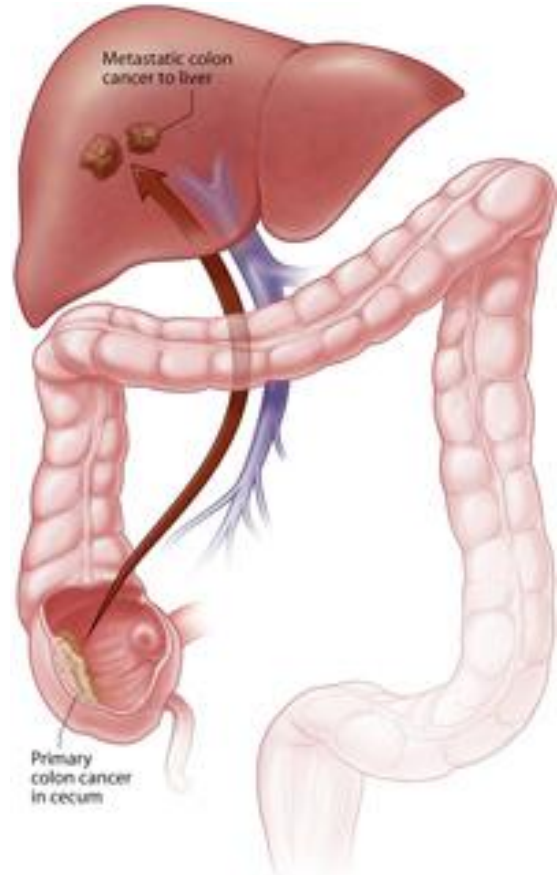
# What makes primary cancer different than metastatic cancer?

Get already processed and summarized RNA-Seq data from *recount2* with **sample information**



Ok. Ok. What about  
actually *using* these  
data and predictions...?

What makes primary cancer different than metastatic cancer?



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



Seon-Kyu Kim<sup>a,1</sup>, Seon-Young Kim<sup>a,1</sup>, Jeong-Hwan Kim<sup>a</sup>, Seon Ae Roh<sup>b,c</sup>,  
Dong-Hyung Cho<sup>c,d</sup>, Yong Sung Kim<sup>a,c,\*\*</sup>, Jin Cheon Kim<sup>b,c,\*</sup>

<sup>a</sup>Medical Genomics Research Centre, Korea Research Institute of Bioscience and Biotechnology, Daejeon, Korea

<sup>b</sup>Department of Surgery, University of Ulsan College of Medicine, Seoul, Korea

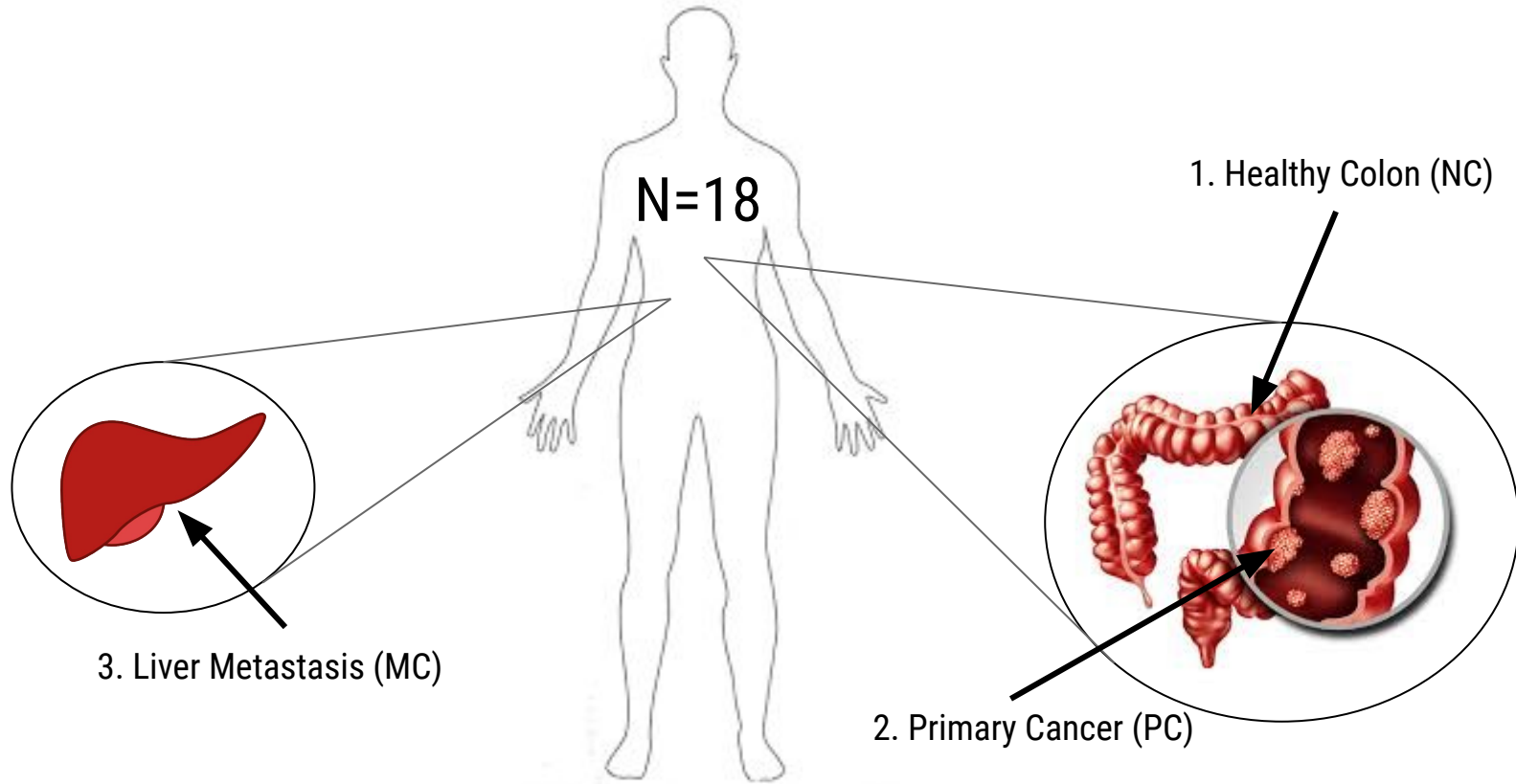
<sup>c</sup>Department of Cancer Research, Institute of Innovative Cancer Research and Asan Institute for Life Sciences, Asan Medical Centre, Seoul, Korea

<sup>d</sup>Graduate School of East-West Medical Science, Kyung Hee University, Gyeonggi-do, Korea

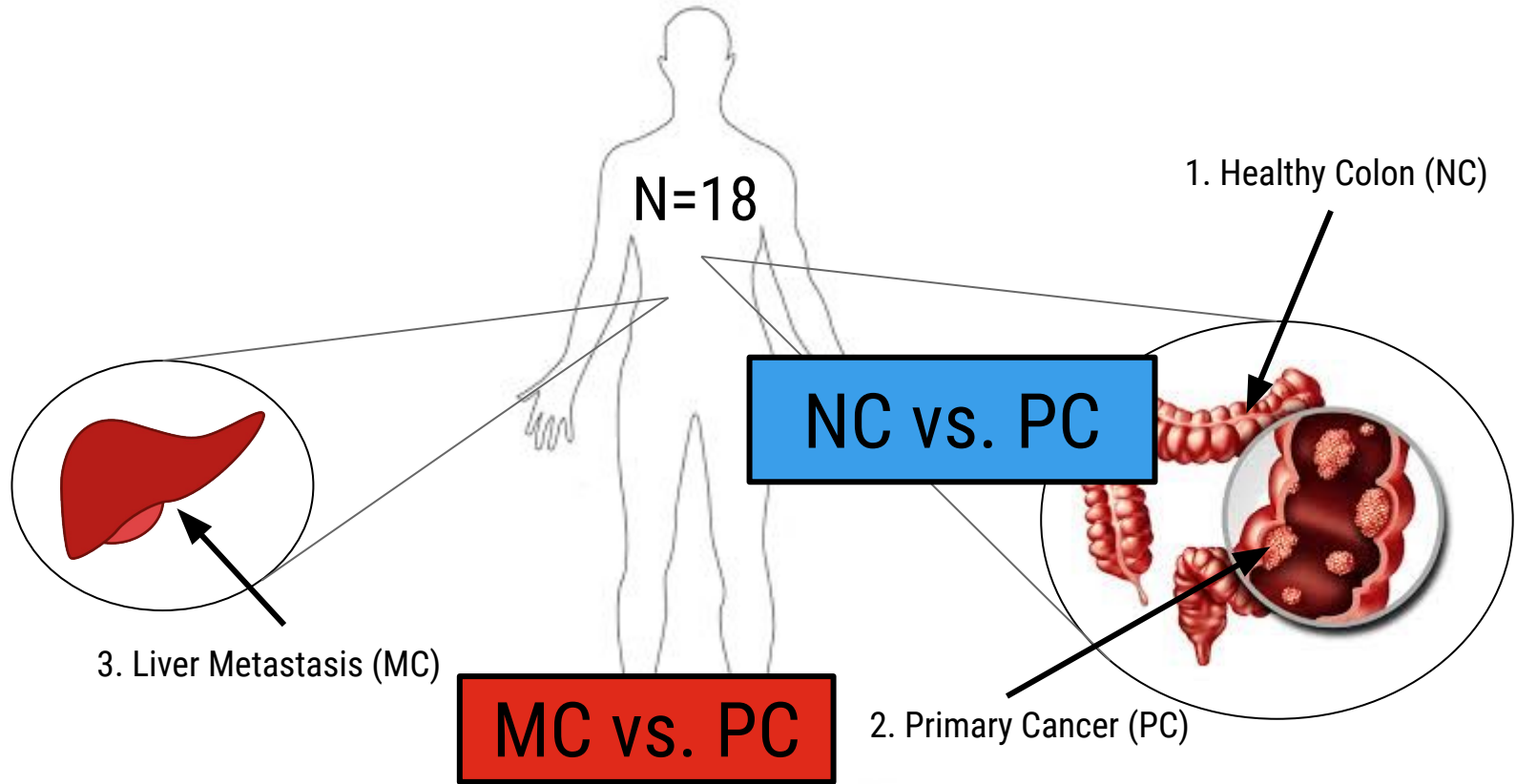
*Molecular Oncology*, July 2014



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



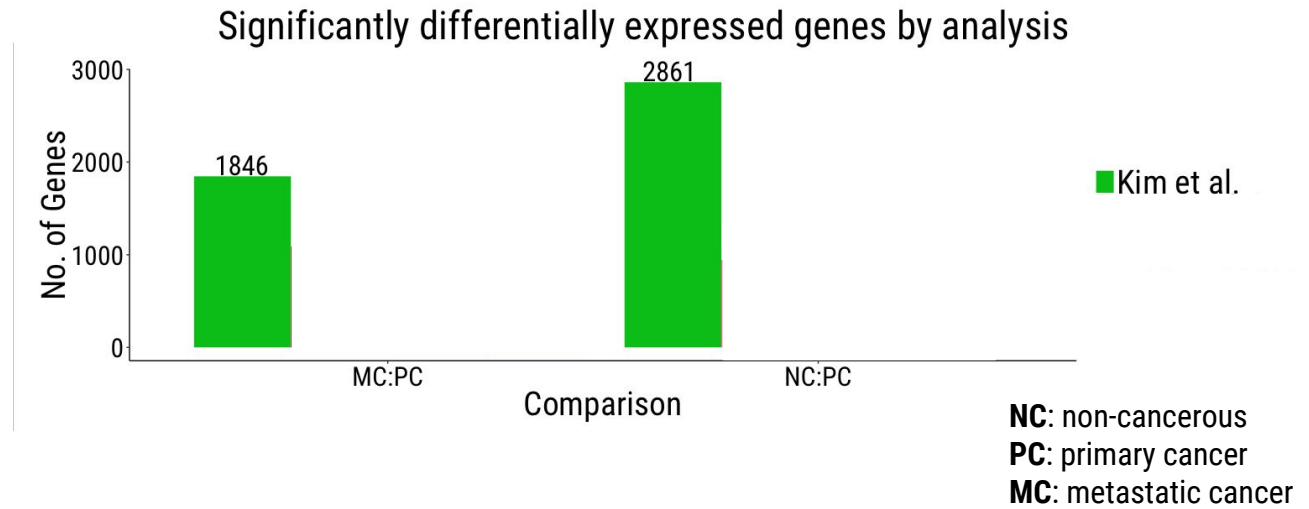
*Kim et al.* analysis looked to identify genes that contribute to metastasis in colon cancer.



Predictions can be used to:

(1) Identify studies of interest

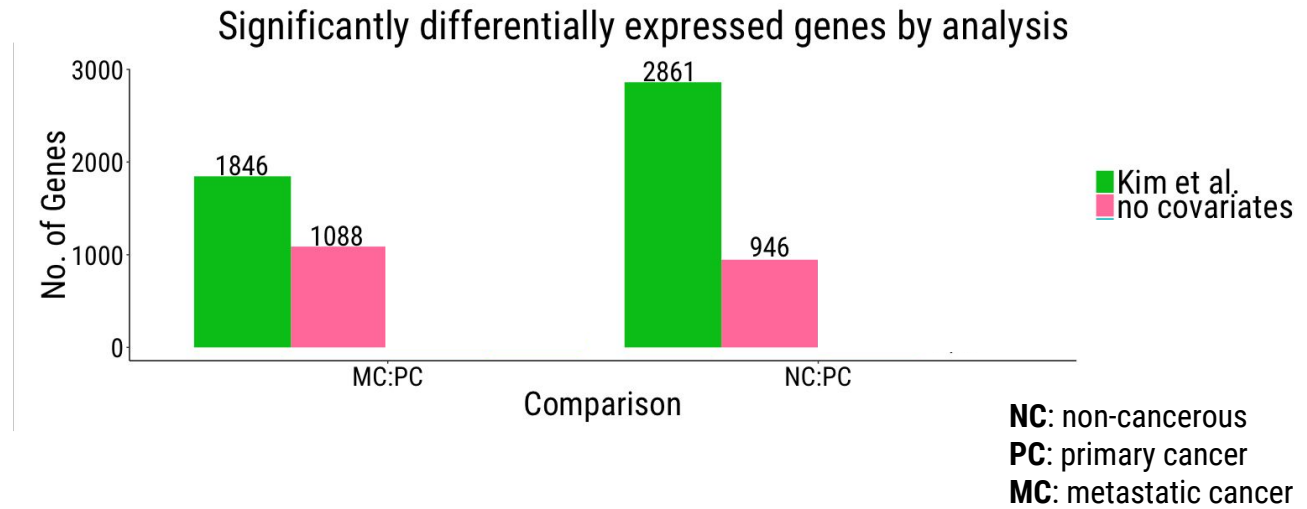
(2) appropriately analyze data



Predictions can be used to:

(1) Identify studies of interest

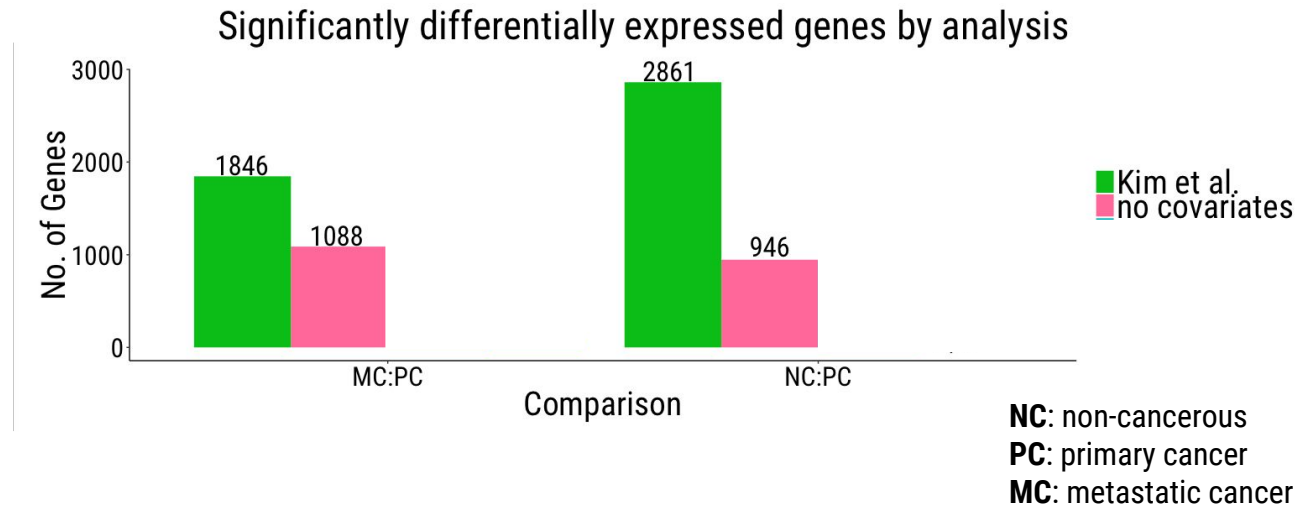
(2) appropriately analyze data



Predictions can be used to:

(1) Identify studies of interest

(2) appropriately analyze data

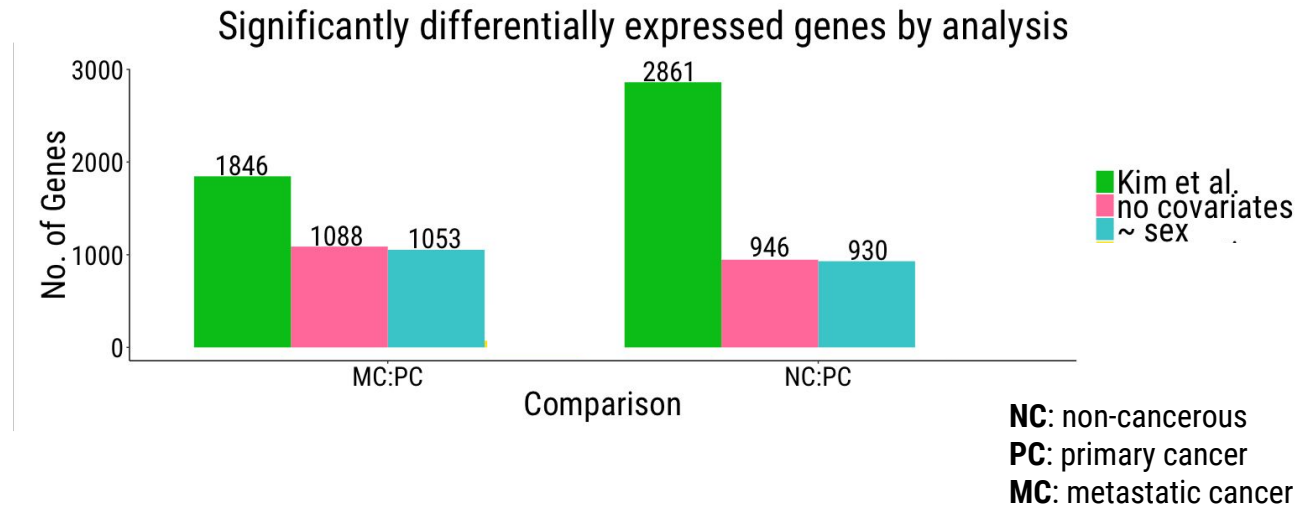


Are the same genes found when sex is included in the analysis?

Predictions can be used to:

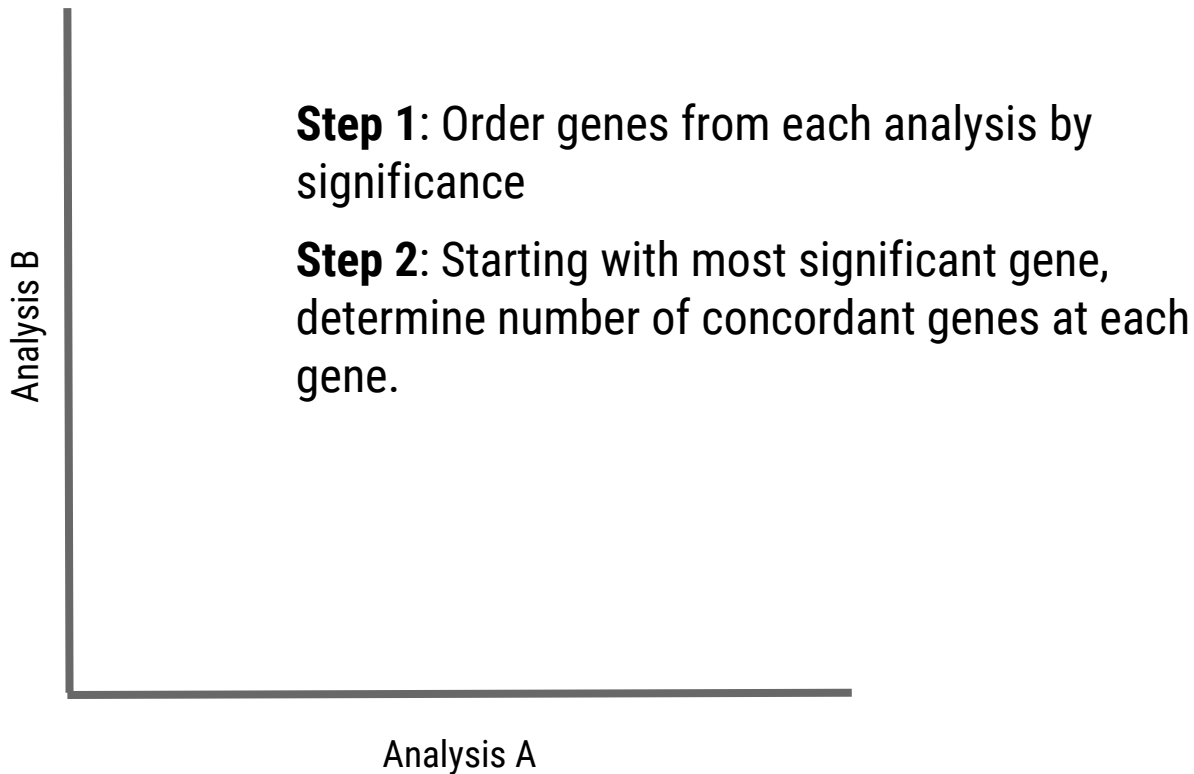
(1) Identify studies of interest

(2) appropriately analyze data



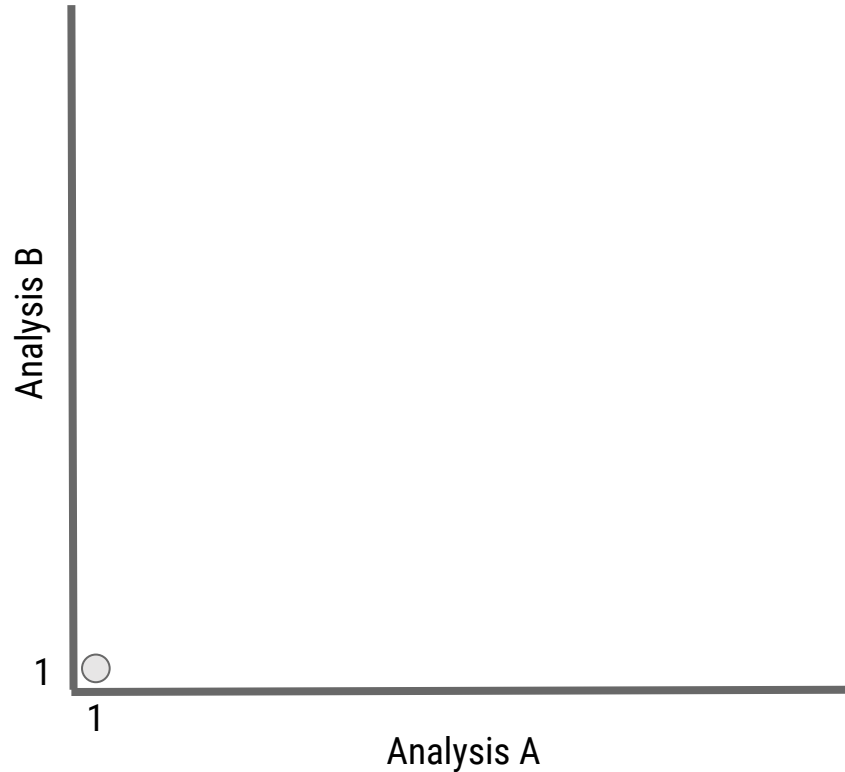
# Concordance at the Top (CAT) Plots

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



# Concordance at the Top (CAT) Plots

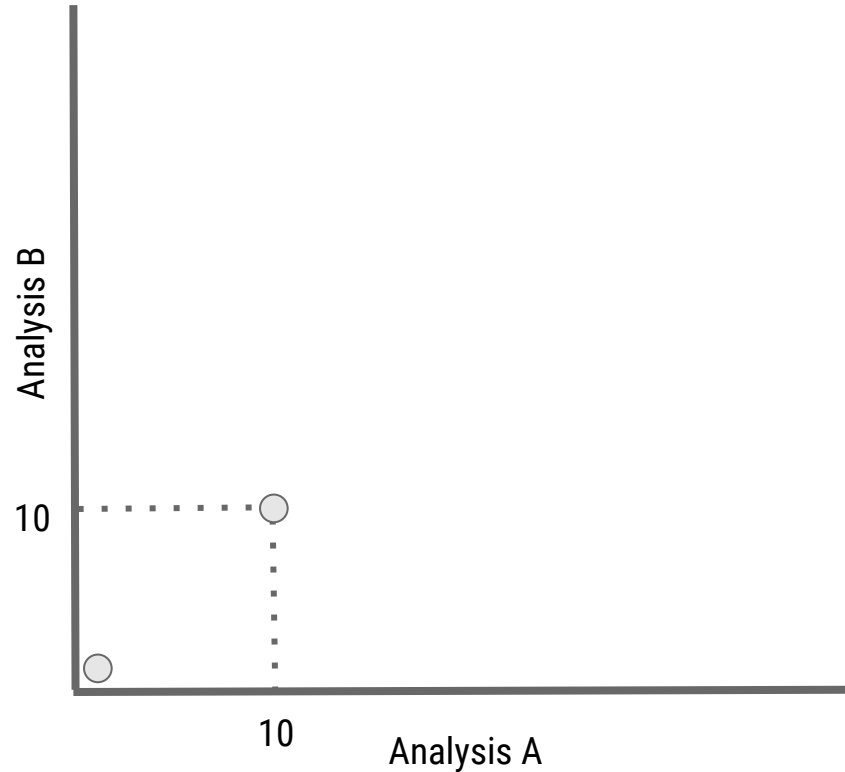
How similar are the results from Analysis A and Analysis B ?





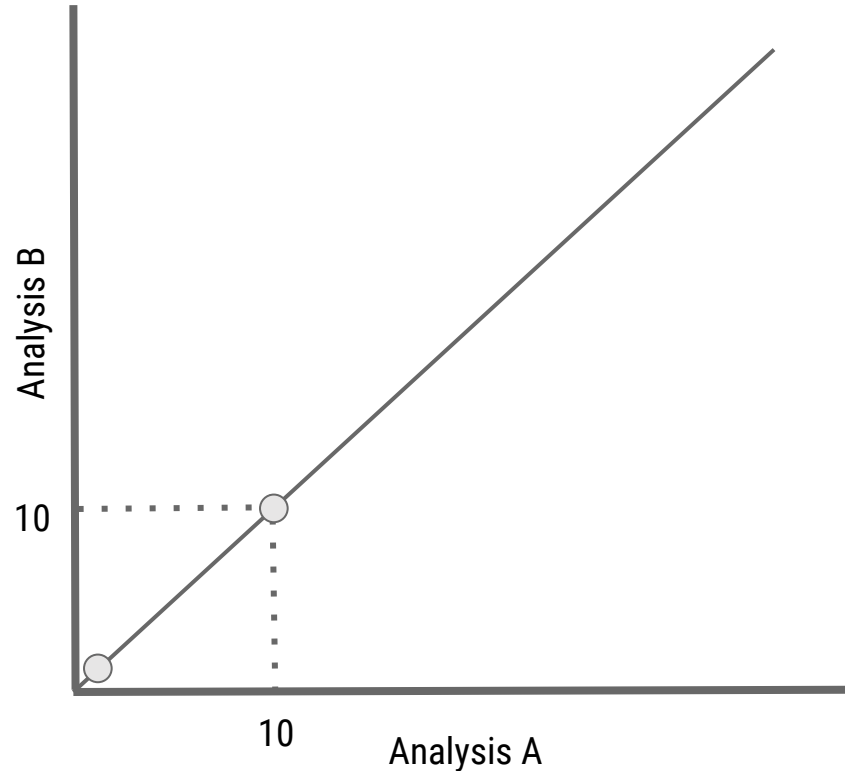
# Concordance at the Top (CAT) Plots

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



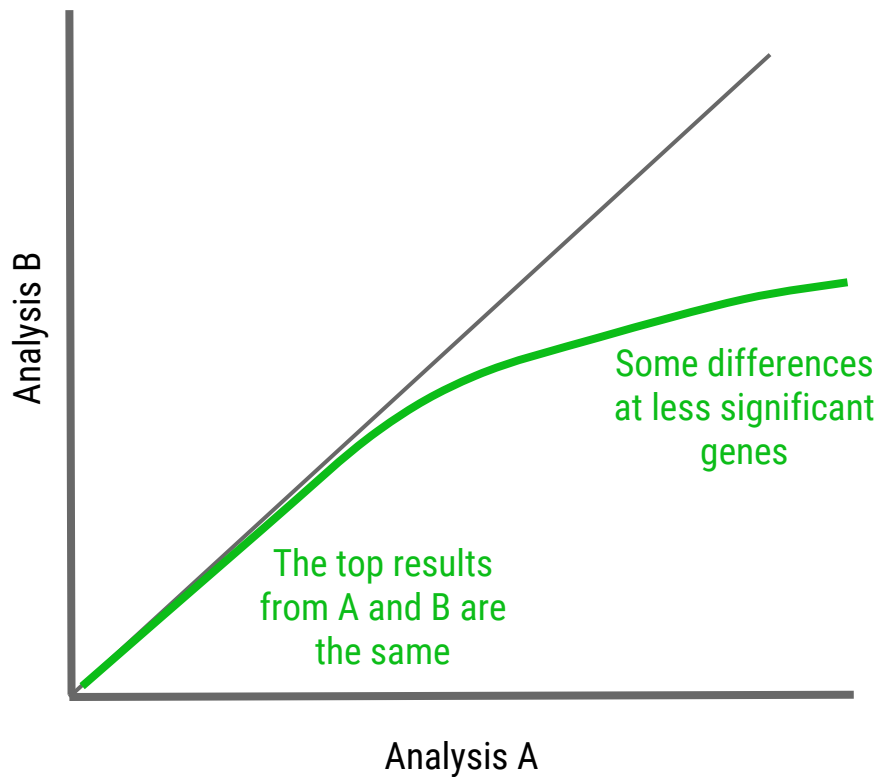
# Concordance at the Top (CAT) Plots

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



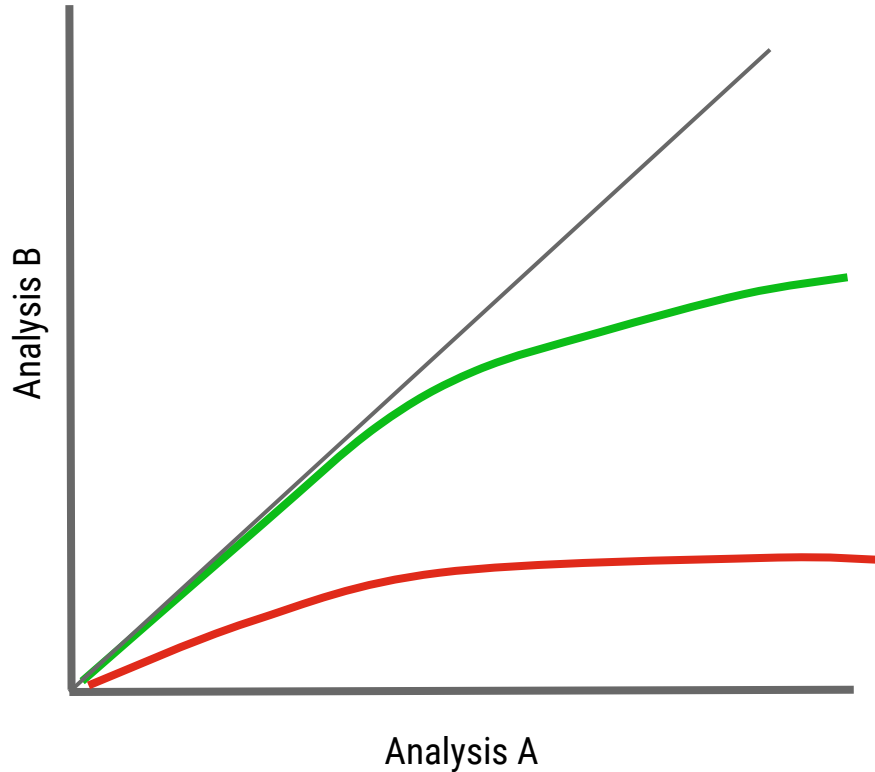
# Concordance at the Top (CAT) Plots

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 ?

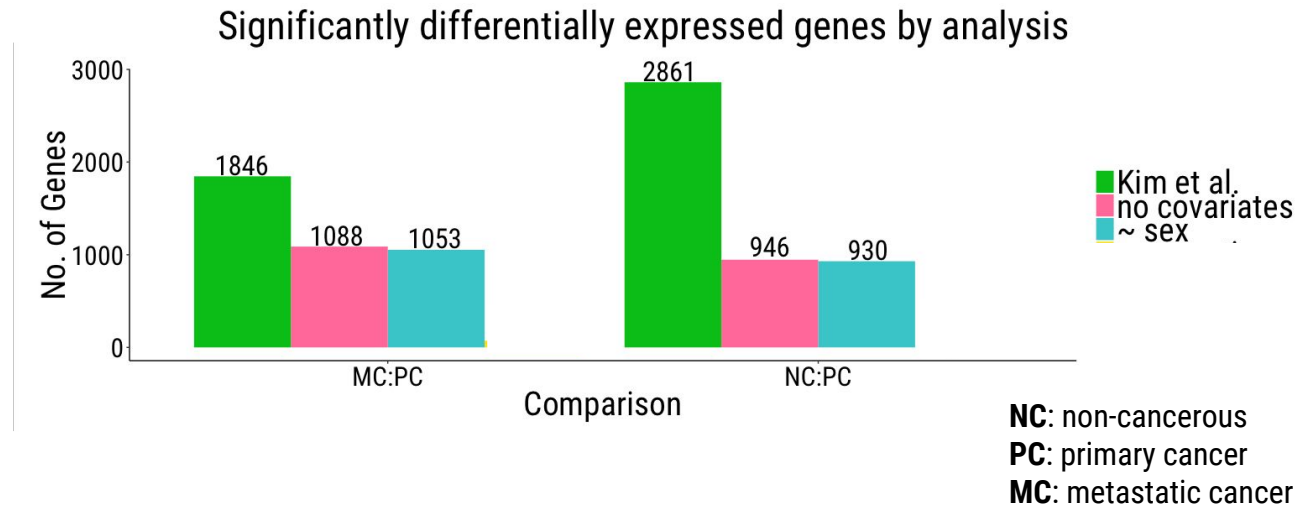


The results of the red condition are less similar between Analysis A and B than the green condition

Predictions can be used to:

(1) Identify studies of interest

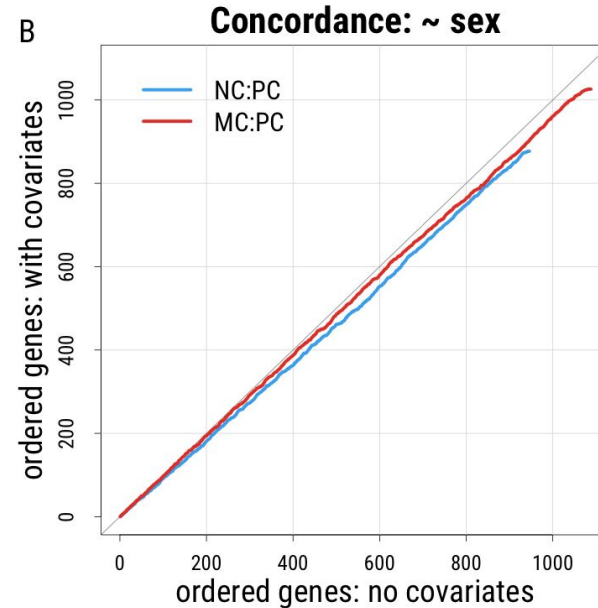
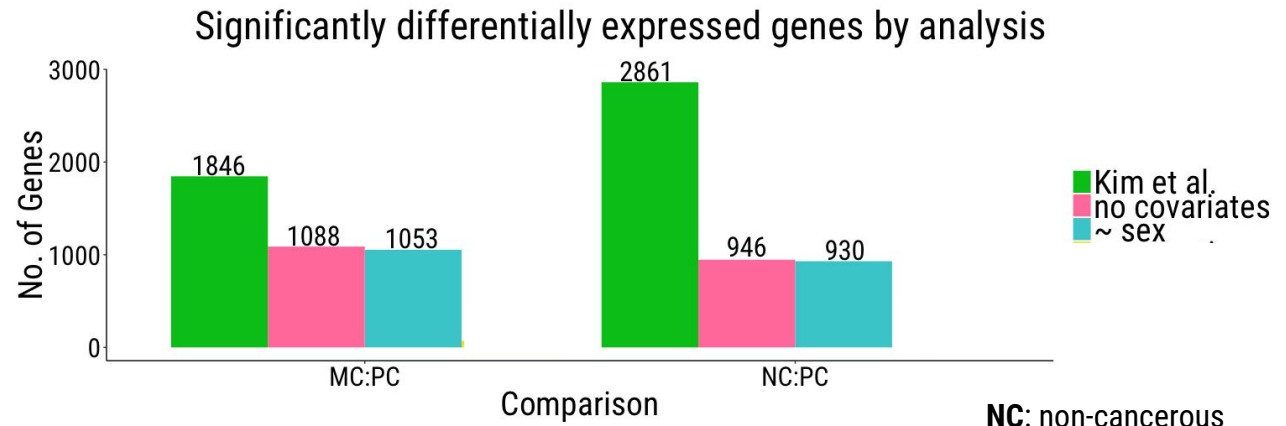
(2) appropriately analyze data



Predictions can be used to:

(1) Identify studies of interest

(2) appropriately analyze data

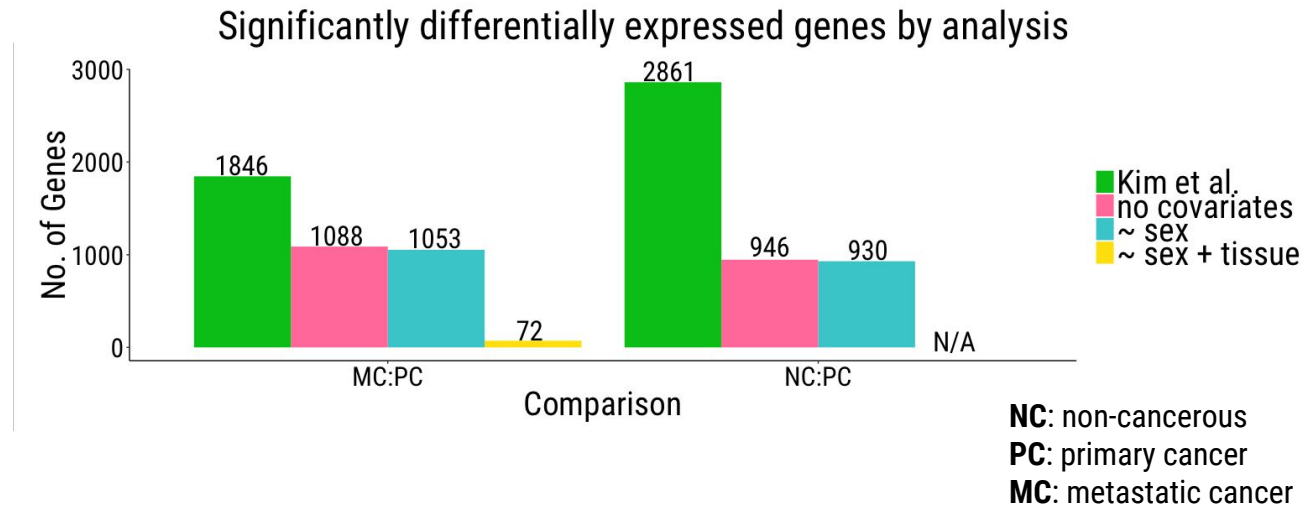


**NC:** non-cancerous  
**PC:** primary cancer  
**MC:** metastatic cancer

Predictions can be used to:

(1) Identify studies of interest

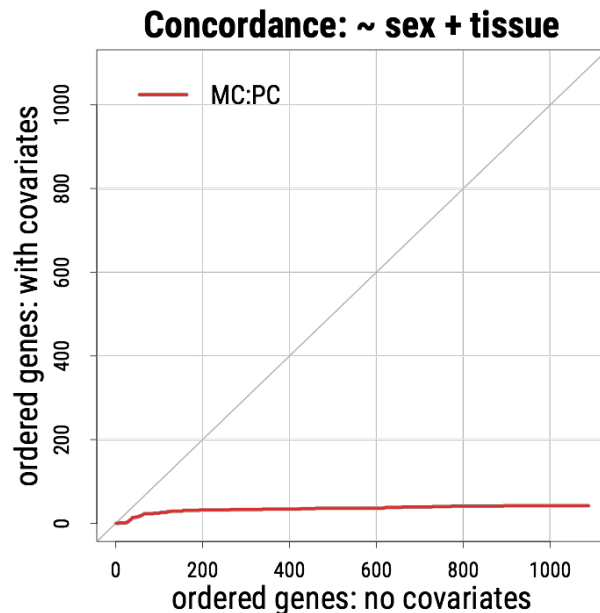
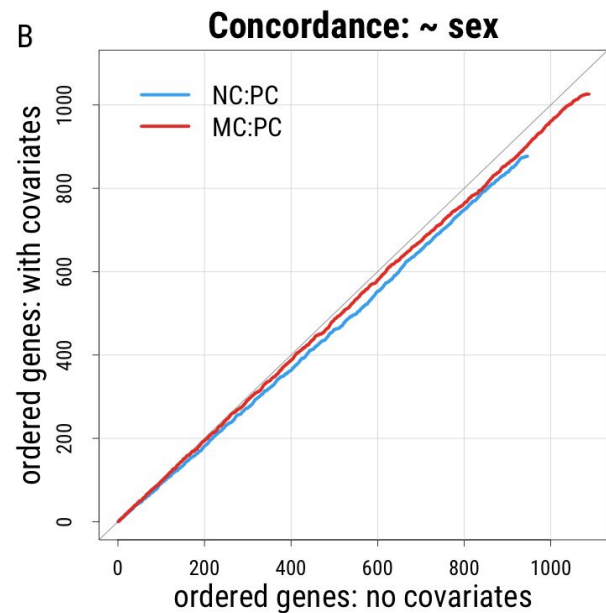
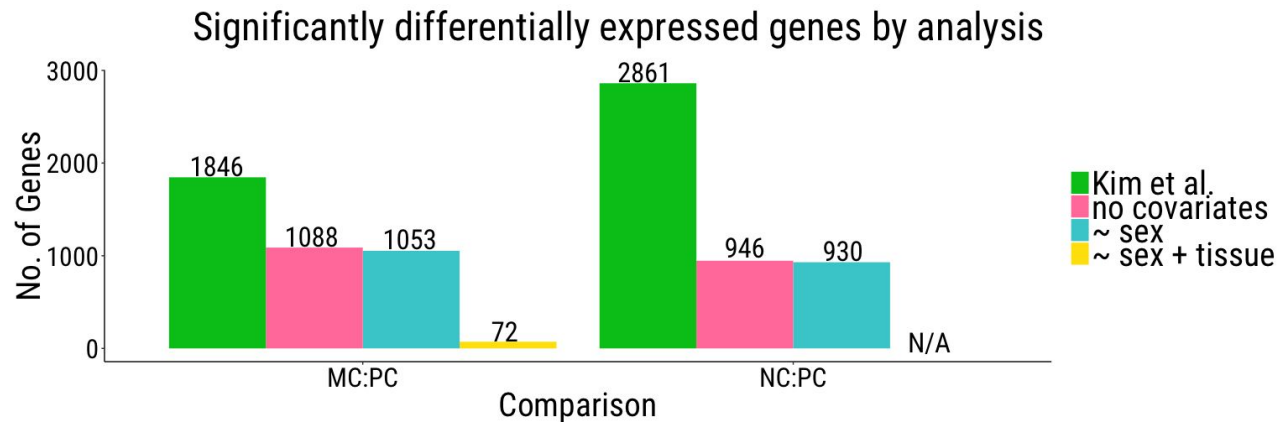
(2) appropriately analyze data



Predictions can be used to:

(1) Identify studies of interest

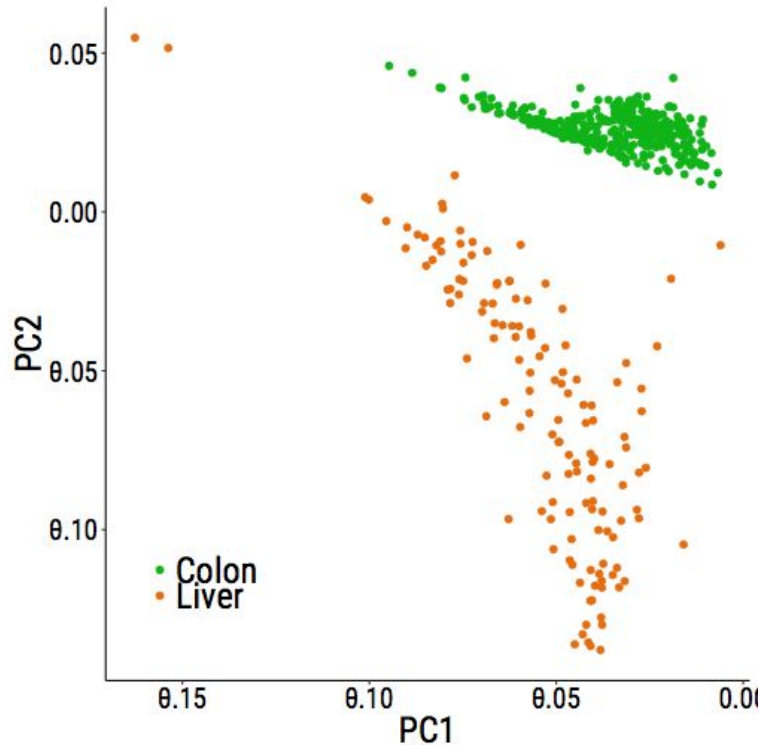
(2) appropriately analyze data



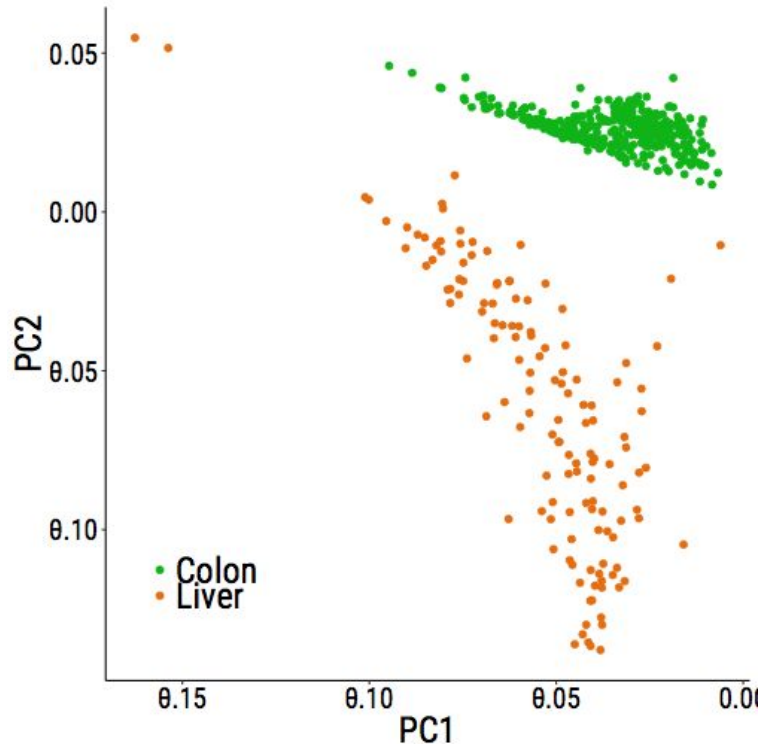


Loss of concordance suggests that differential expression is detecting tissue differences, not cancer-related changes.

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

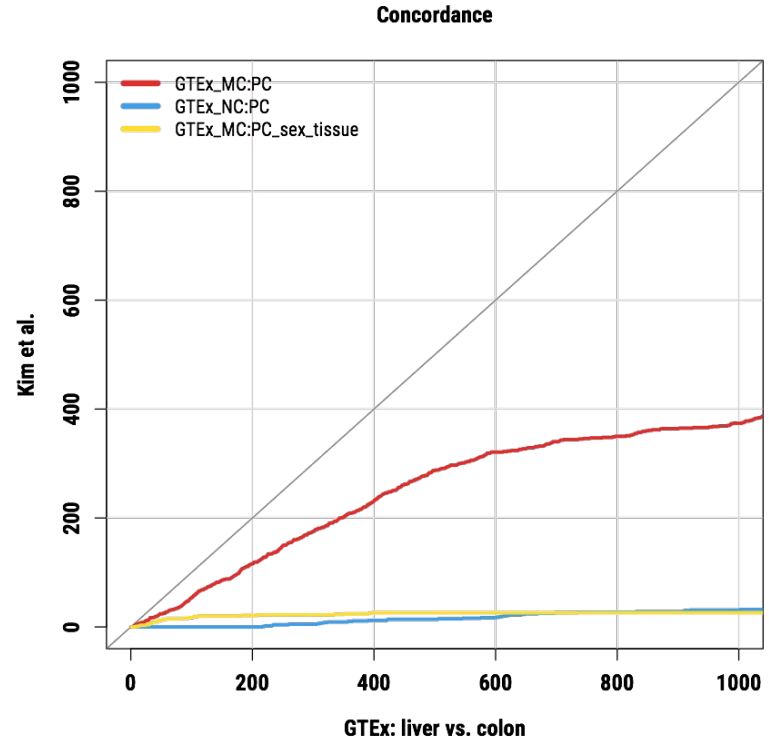
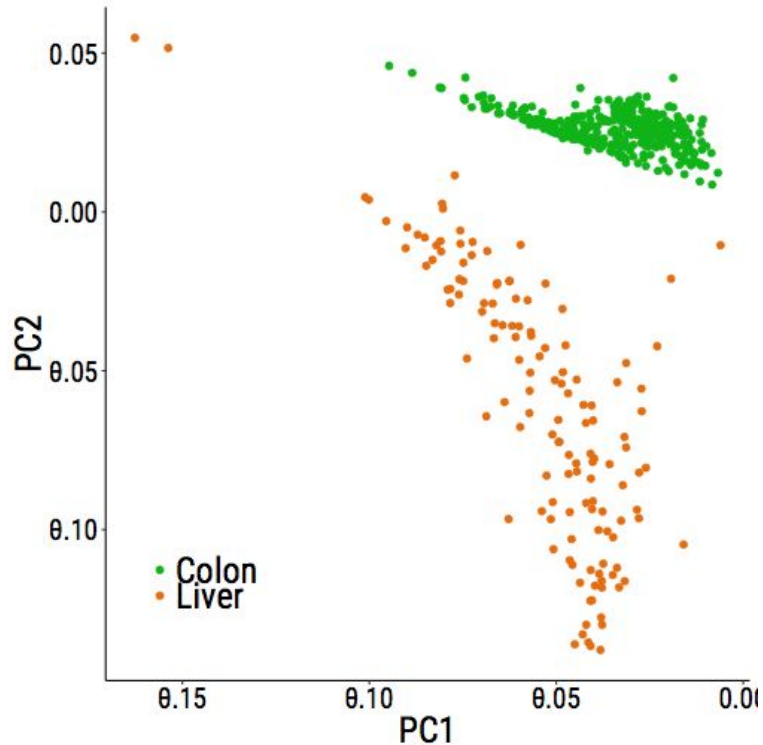


So...what if we compared the MC:PC results with differential expression between colon and liver?



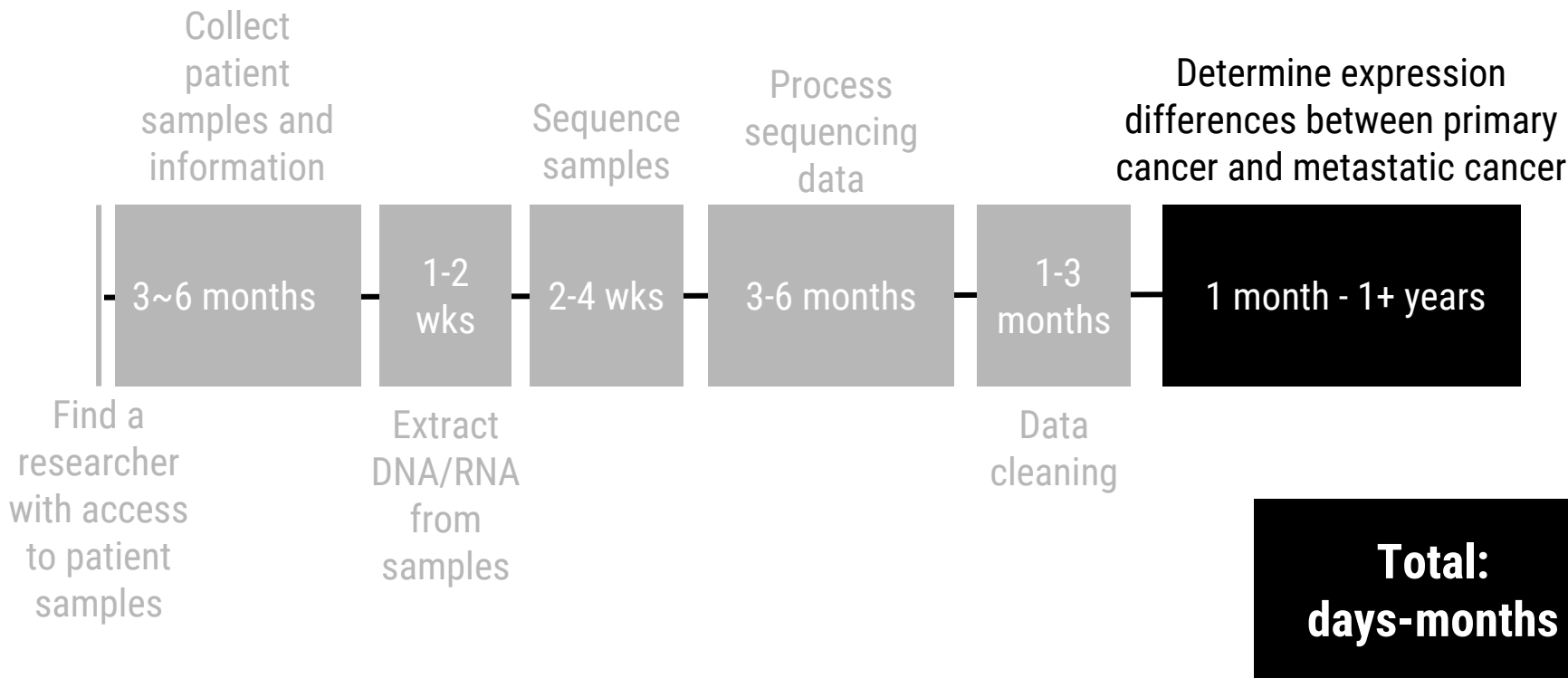
**Hypothesis:** MC:PC results should be most similar to GTEx colon vs. liver

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





# What makes primary cancer different than metastatic cancer?



## The Leek group

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

## Collaborators

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

# A quick tour of a geneticist turned data scientist

## Background

## Projects

1. PhD work studying the genetic basis of autism
2. Postdoctoral work working with 70,000 samples
3. Working toward accessible data science education

## What I do here at UCSD



# Chromebook Data Science (CBDS)





Find a **partner organization**



Collaboratively **develop course content**



Develop **new technology** as needed



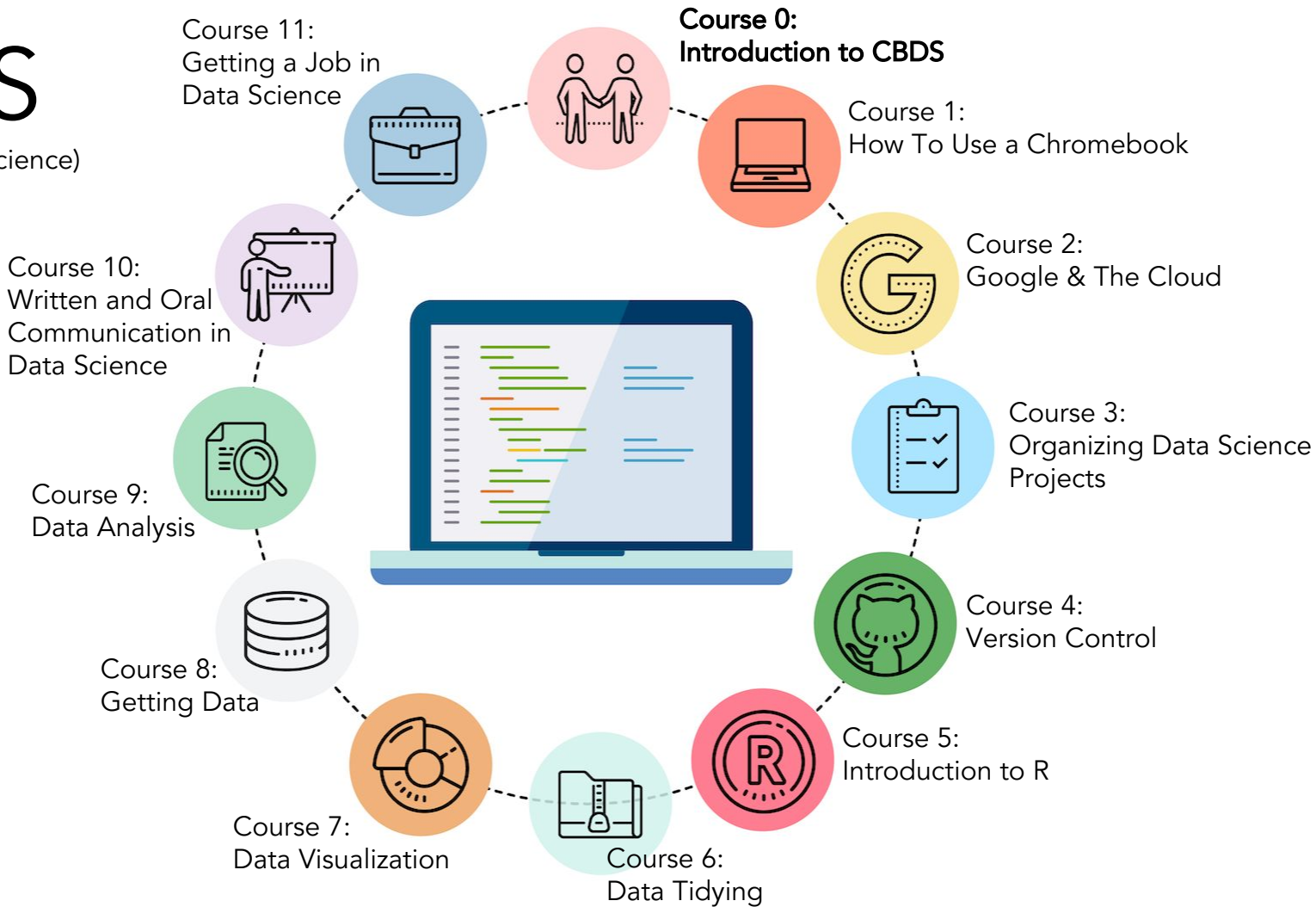
Design in-person **tutoring program**



Launch program, **teach the stuff** & get learners **jobs**

# CBDS

(Chromebook Data Science)



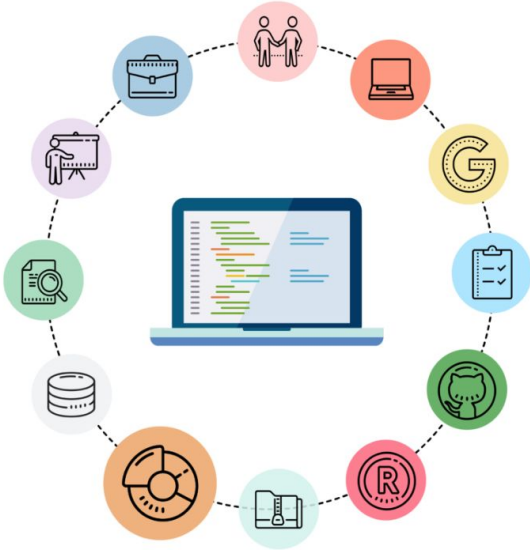


Library

Author

Community

Support



# Chromebook Data Science

## Course 7: Data Visualization

### Data Visualization



Jeffrey Leek

\$0.00 / MINIMUM

\$25.00 / SUGGESTED ⓘ

YOU PAY

\$25.00

YOU PAY (US\$)

\$25.00

EU customers: Price excludes VAT.  
VAT is added during checkout.

Add Course To Cart

The instructor has published 100% of this course.

Course Info

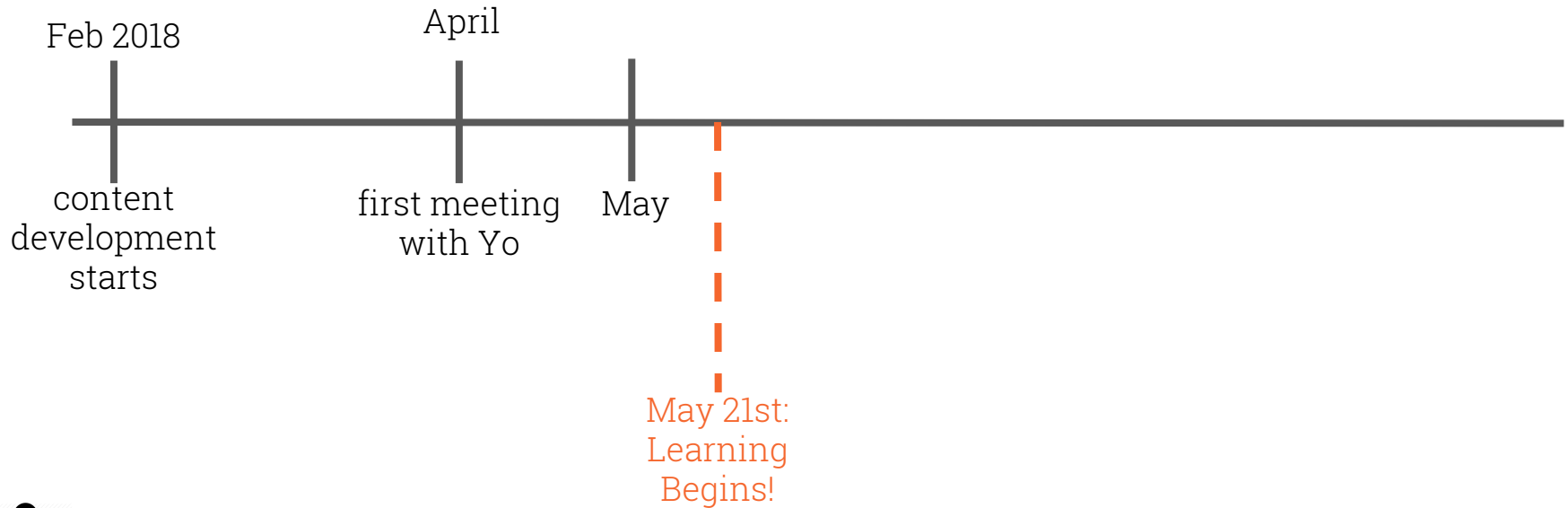
Course Materials

Instructors

Community

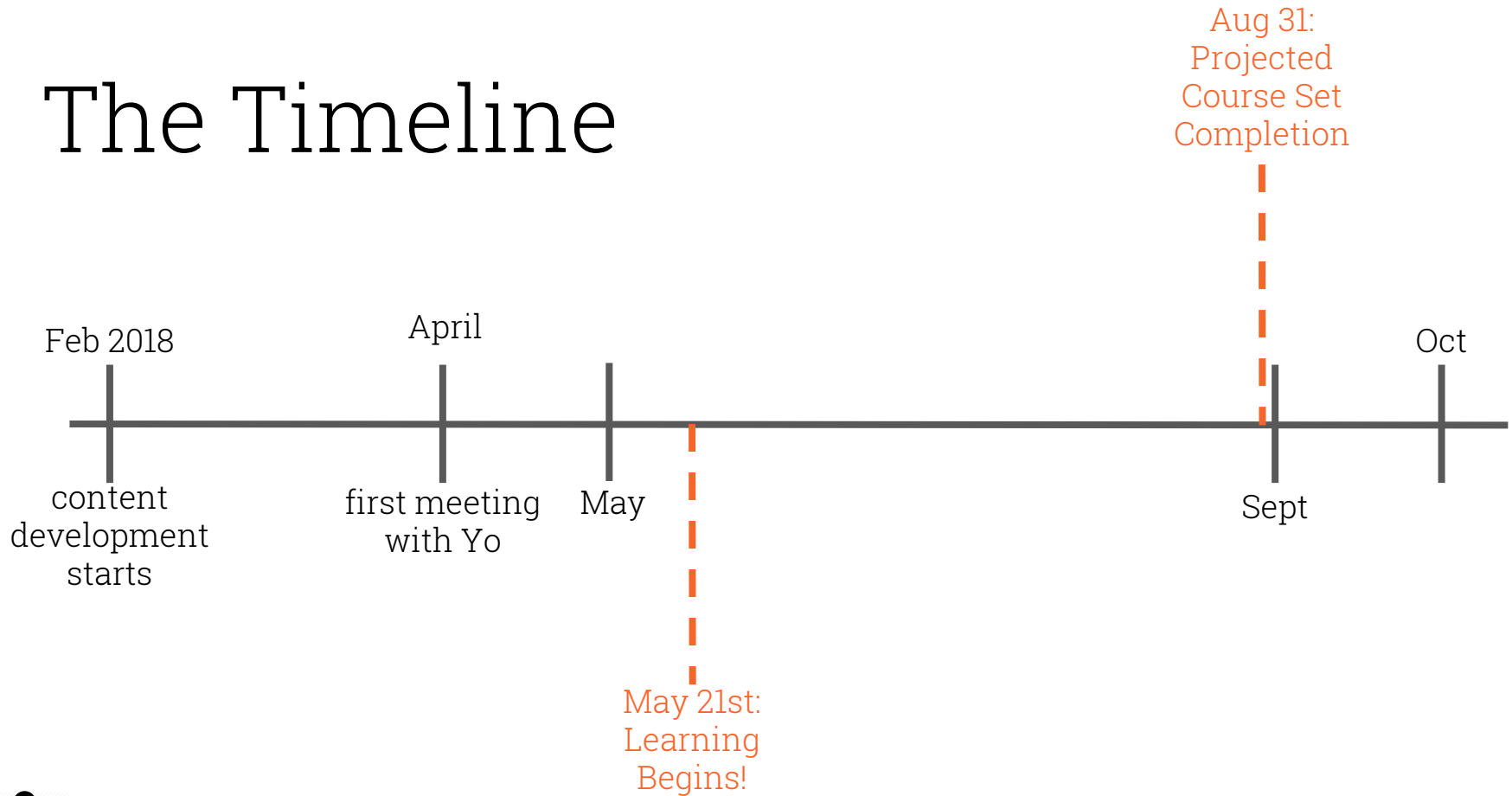


# The Timeline



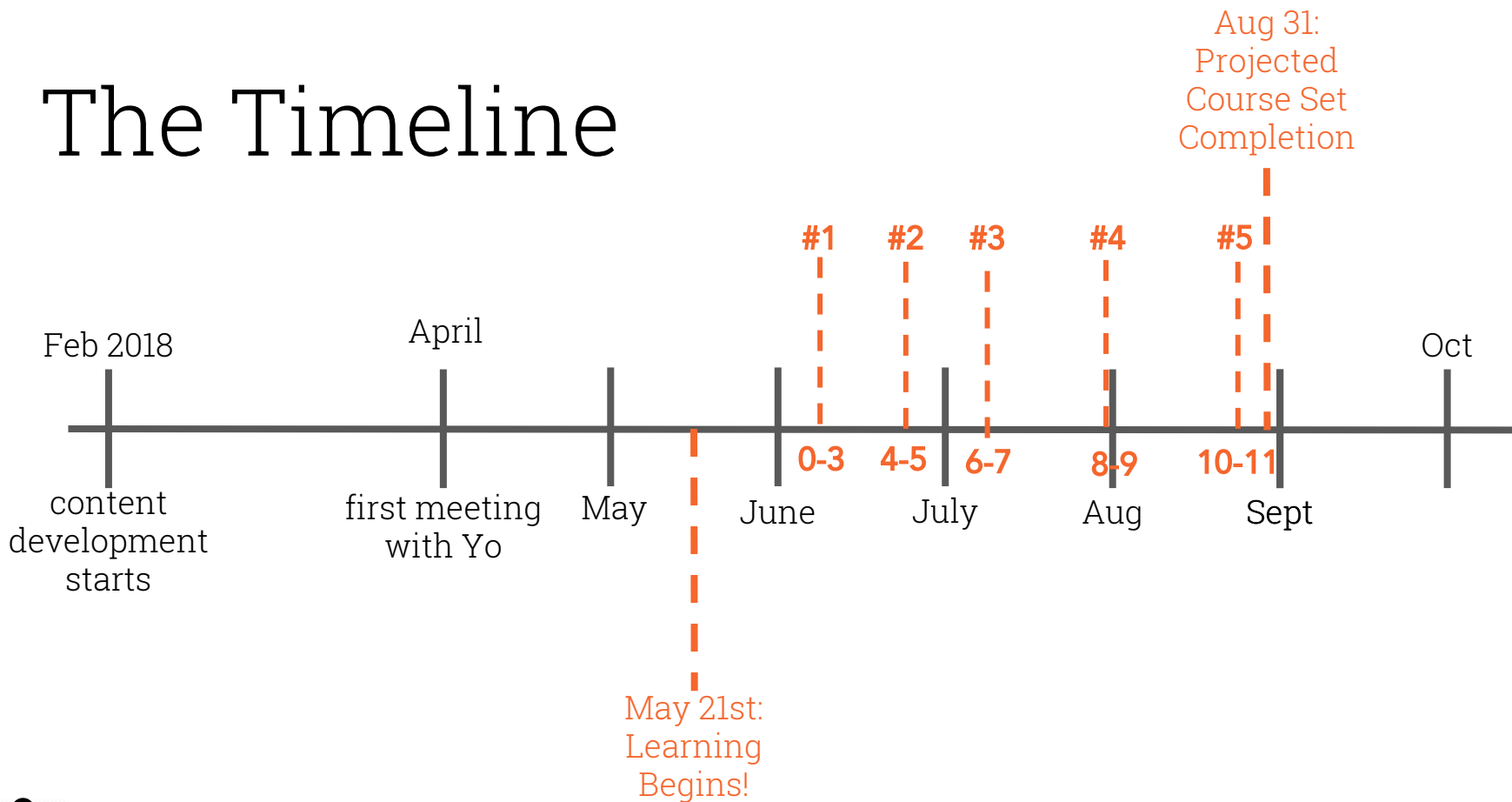
Launch program, **teach the stuff** & get learners **jobs**

# The Timeline



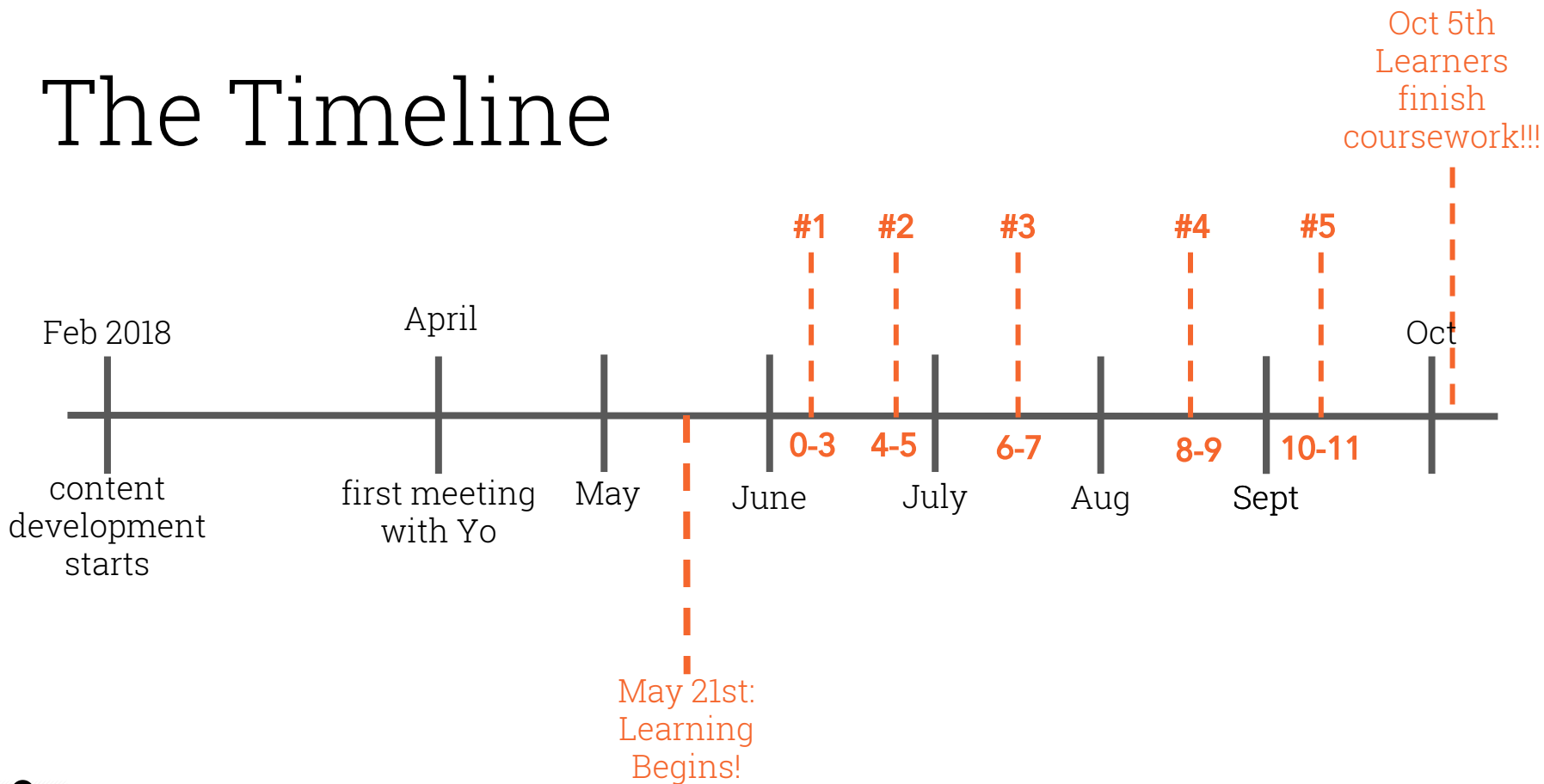
Launch program, **teach the stuff** & get learners **jobs**

# The Timeline



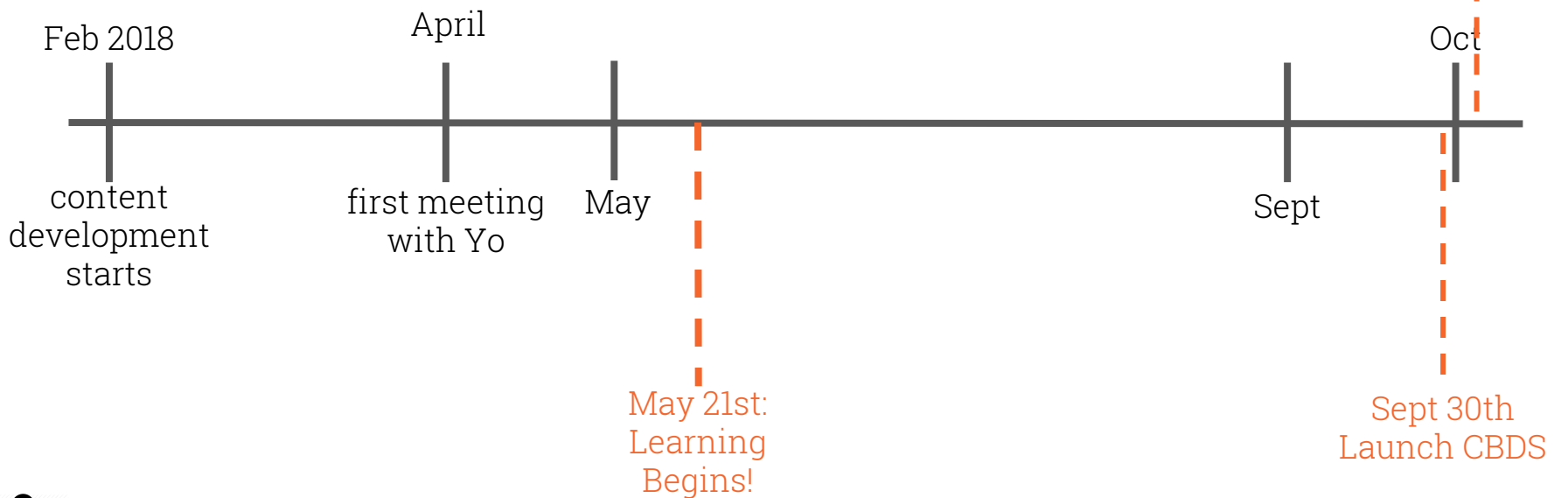
Launch program, **teach the stuff** & get learners **jobs**

# The Timeline



Launch program, **teach the stuff** & get learners **jobs**

# The Timeline



Launch program, **teach the stuff** & get learners **jobs**





john "as b/4" fink

@adr

Follow



rstudio and Chromebook Data Science are so fun I am working on them on a Saturday night so thanks a buncha buncha @kierisi for the recommendation.

9:40 PM - 6 Apr 2019

2 Retweets 10 Likes



# Chromebook Data Science Team



Leek

Johnson

Ellis

Hadavand

Muschelli

Kross

Collado-Torres

Jager

McClymont

Myint

Content Development



Administration & Tutoring



Technology



# A quick tour of a geneticist turned data scientist

## Background

## Projects

1. PhD work studying the genetic basis of autism
2. Postdoctoral work working with 70,000 samples
3. Working toward accessible data science education

## What I do here at UCSD

# COGS9 : Intro to Data Science

Brad has previously written about [COGS9](#) in his super interesting and thoughtful blog post [Data Science at UC San Diego](#), but very briefly here, COGS9 is an intro level course designed to get undergraduate students interested in data science, familiar with what data science is, and excited to learn more. It is neither math nor computationally-heavy, but is rather taught through concepts and examples. When it first ran there were 24 students. Now, each time it is offered, there are hundreds. This quarter, when numbers settled out, I had 326 students in the course.



# Welcome to COGS18: Introduction to Python!

COGS 18 · Shannon Ellis · Spring 2019 · UCSD

COGS 18 - Introduction  
To Python

[Home](#)

[Materials](#)

[Coding Labs](#)

[Assignments](#)

## Overview

Introduction to Python (COGS18) is a course offered by the Department of Cognitive Science of UC San Diego, taught by [Shannon Ellis](#). It is a hands-on programming course, focused on teaching students in Cognitive Science and related disciplines an introduction on how to productively use Python.

## Current Iteration

Introduction to Python is currently running for Spring Quarter 2019, for which you can check out the current [syllabus](#) and [schedule](#). Course lectures are recorded and are publicly available as screencasts from [here](#).



# COGS108 - Data Science in Practice

Course materials for Hands-On Data Science.

UC San Diego    COGS108@gmail.com

Repositories 23

People 19

Teams 4

Projects 0

Settings

## Pinned repositories

Customize pinned repositories

### Overview

Overview and map of the organization, which services COGS108: Hands-On Data Science, from UCSD.

★ 41    🍴 17

### Lectures-Sp19

Slides and Notebooks used in Lecture for Sp19 COGS108

★ 5    🍴 22

### Section\_Workbooks

Workbooks for practice during discussion section

● Jupyter Notebook    🍴 25

### Tutorials

Tutorial notebooks for hands-on data science, following along with the course topics.

● Jupyter Notebook    ★ 38    🍴 117

### Projects

Final Project materials and description.

● Jupyter Notebook    ★ 3    🍴 95

### Readings

A curated list of suggested reading materials.

★ 6    🍴 2

# UCSD



**Brad Voytek**



**Tom Donoghue**



**Instructional staffs, students,  
& COGS faculty/staff**