

Data Centric View of DNA/RNA Information Flow

F. Alex Feltus, Ph.D.

CU-HackIt 2022

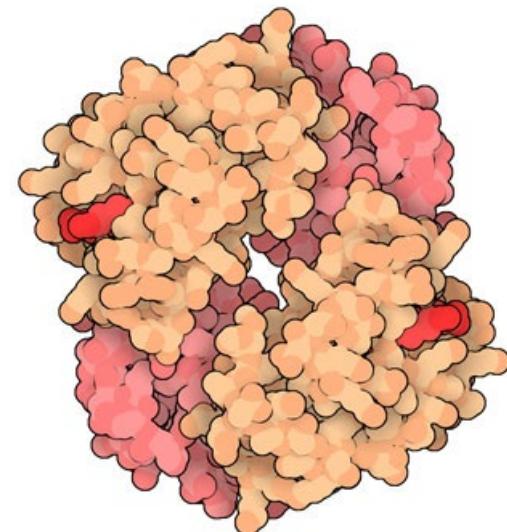
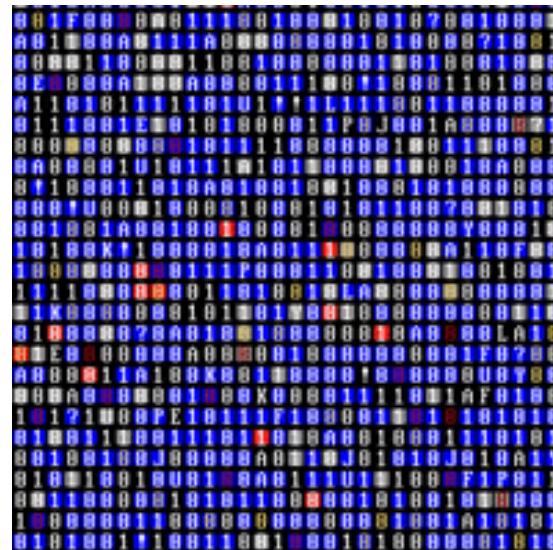
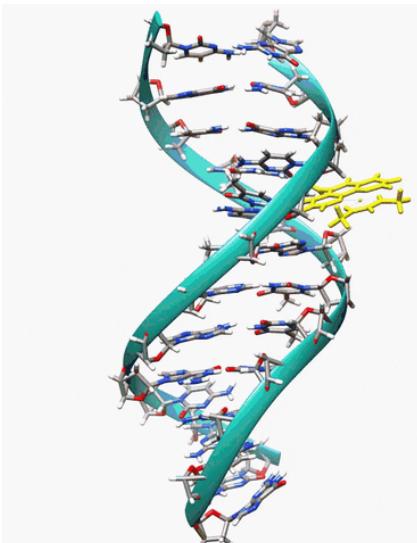
Professor, Clemson Dept. of Genetics & Biochemistry

Co-Founder Praxis AI LLC

ffeltus@clemson.edu

Important point:

The cellular “hard drive” is molecular. Molecular storage devices like DNA & RNA are sequences (aka strings) that are perfect for analysis on computers!



A close-up photograph of a fluffy, blue and white cat lying on its side. The cat's eyes are closed, and it appears to be sleeping or resting. Its soft fur is visible, and it looks very peaceful.

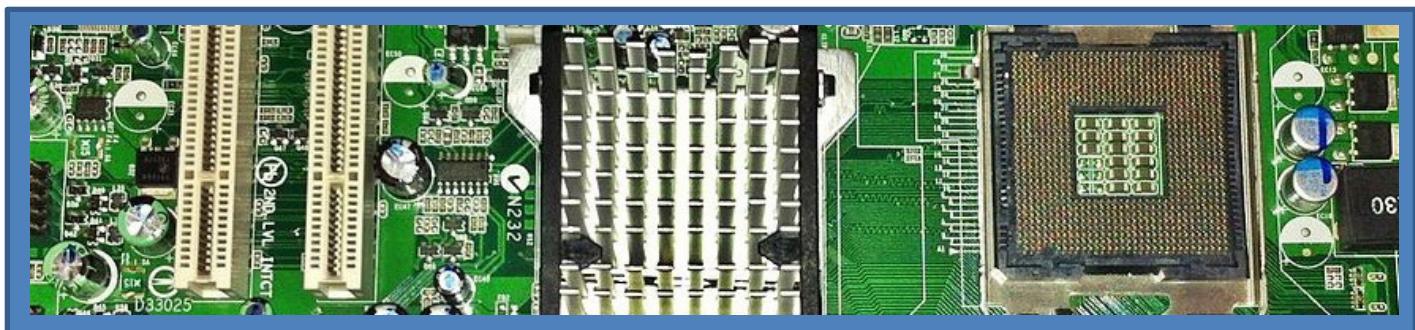
Who is this
guy and why
should I not
be bored?

I Run A **Small** Interdisciplinary Computational Biology Research Lab

Mammals

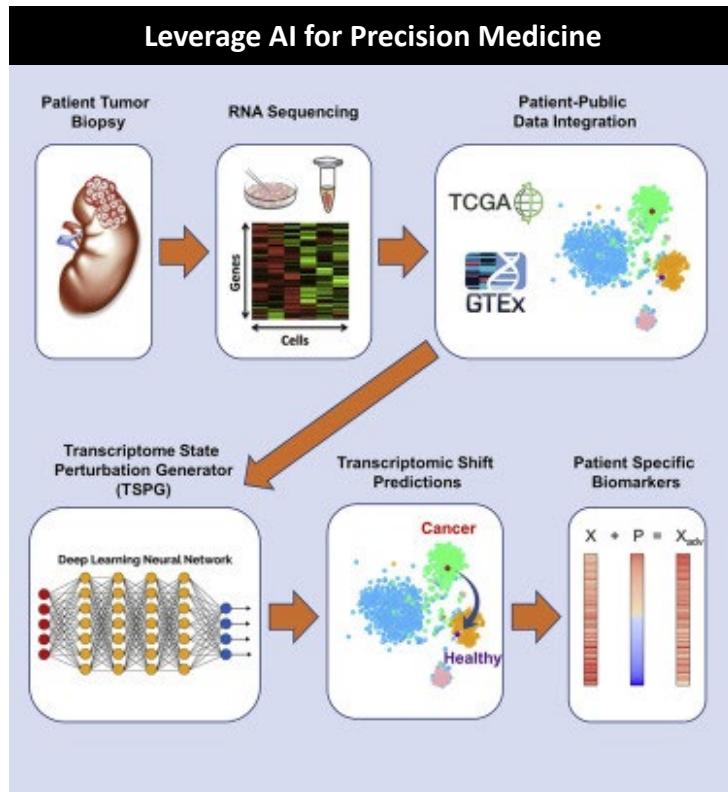


Plants



Bioinformatics + Cyberinfrastructure Engineering

Example Interdisciplinary Research in Feltus Lab @ Clemson



Patterns

Available online 17 August 2020, 100087

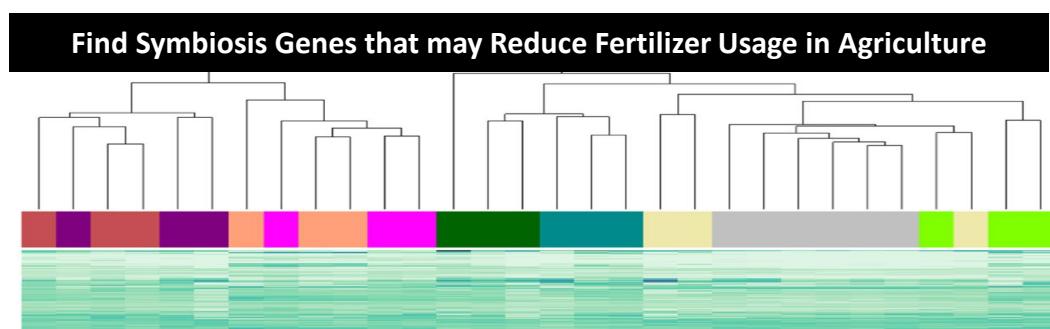
In Press, Corrected Proof



Article

Cellular State Transformations Using Deep Learning for Precision Medicine Applications

Colin Targonski ^{1, 4}, M. Reed Bender ^{2, 6}, Benjamin T. Shealy ¹, Benafsh Husain ², Bill Paseman ³, Melissa C. Smith ¹, F. Alex Feltus ^{2, 4, 5, 7, 8}



Identifying Temporally Regulated Root Nodulation Biomarkers Using Time Series Gene Co-Expression Network Analysis

William L. Poehlman[†], Elise L. Schnabel, Suchitra A. Chavan, Julia A. Frugoli^{*} and Frank Alex Feltus[†]

Scalable Computational Biology Training

Clemson
11-12th Grade
HBCUs
TCUs
<R1 SC Campuses

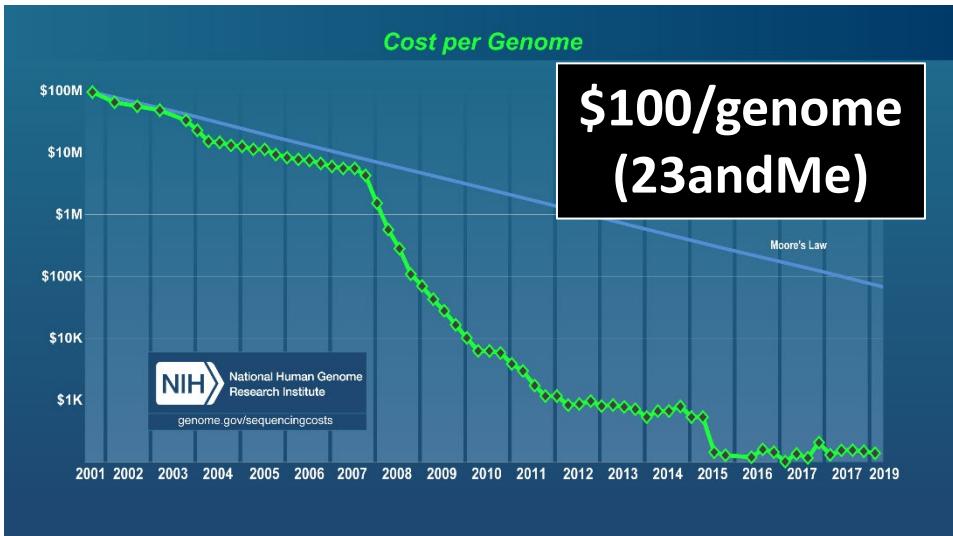
PRAXIS

prxai.com

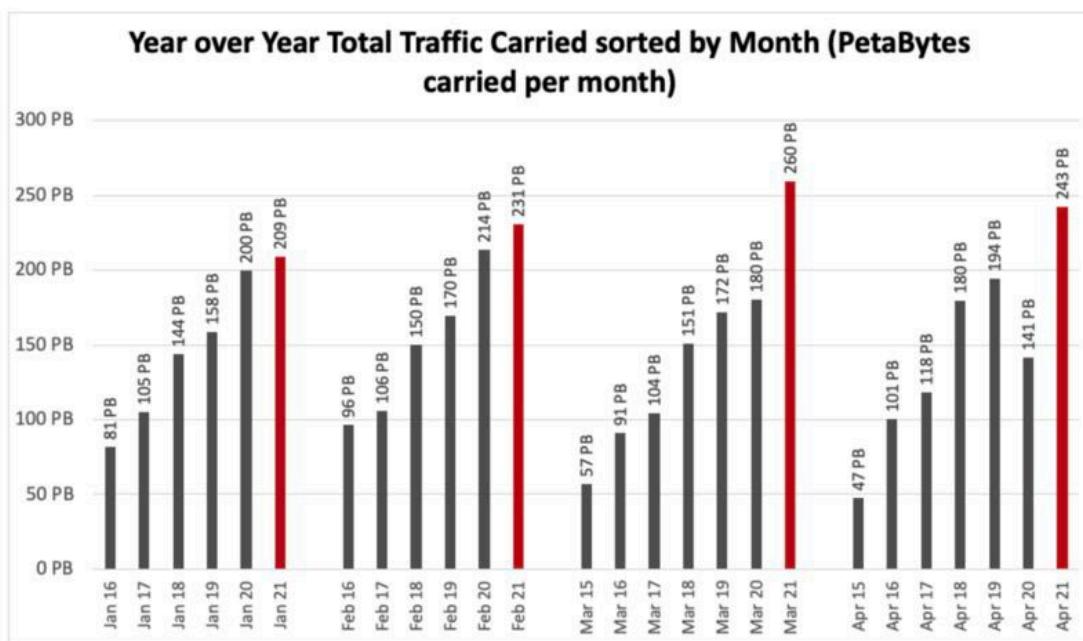
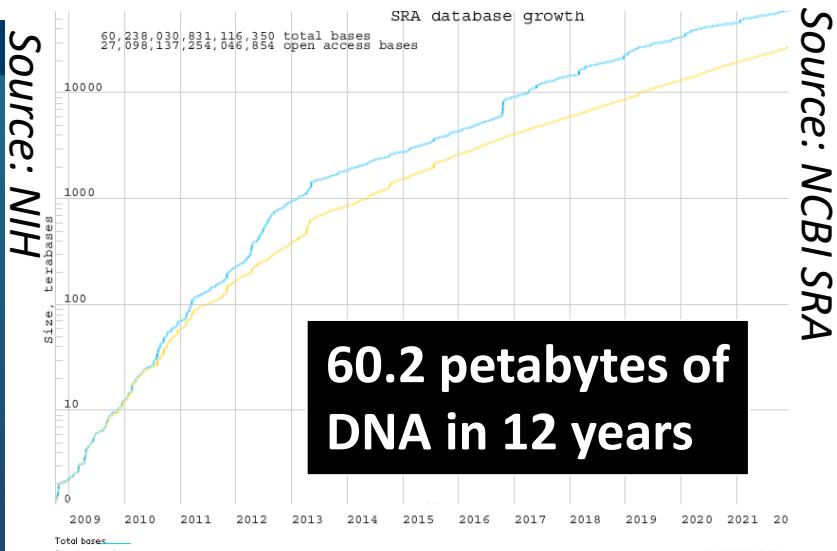
+Brain and Cancer Genomics!

A Lot Has Happened in Biology in the Last Decade

DNA Sequencing Cost is Dropping



DNA Databases are Swelling



Internet2 Traffic > 1 Exabyte Jan-May 2021

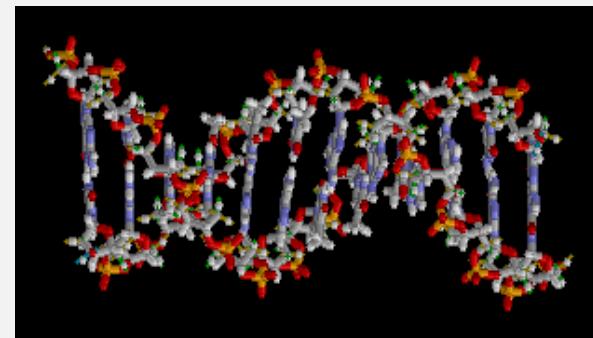
Source: NCBI/SRA

Bioinformatics Requires Advanced Cyberinfrastructure

Cyberinfrastructure is a research environment that supports advanced data acquisition, data storage, data management, data integration, data mining, data visualization and other computing and information processing services distributed over the Internet beyond the scope of a single institution. -wikipedia



Bioinformatics is the technology that uses computers for storage, retrieval, manipulation, and distribution of information related to biological macromolecules (DNA, RNA, proteins, metabolites, and lipids).



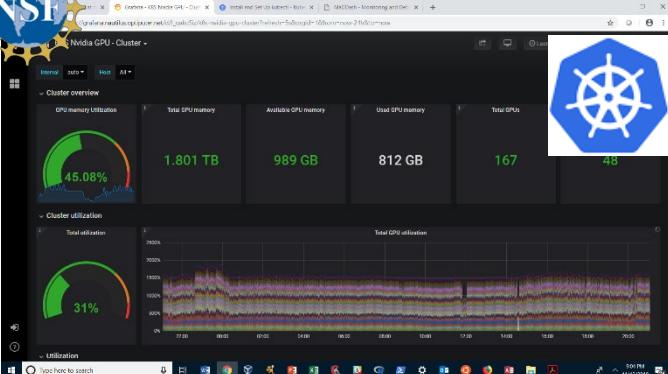
We use the Cloud and Supercomputers to do our work



Clemson Palmetto Cluster



PRP/NRP Kubernetes Cluster



Open Science Grid



HT CENTER FOR
HIGH THROUGHPUT
COMPUTING

chtc.cs.wisc.edu

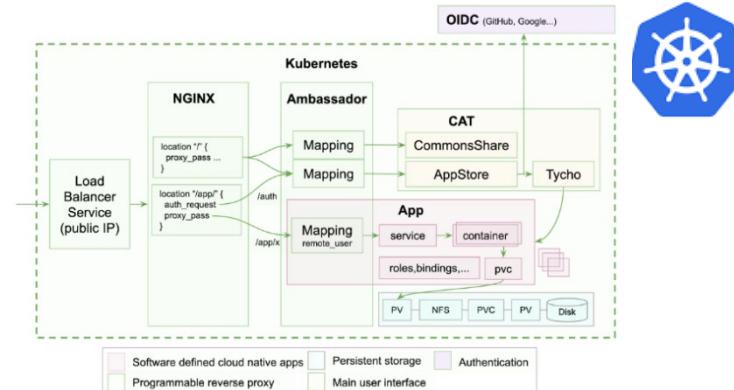


In 2017 on OSG ...

8.43 Million Wall Hours
4.50 Million CPU Hours
8.92 Million Jobs
16.6 Million Transfers
4.07 PB



SciDAS (Scientific Data Analysis at Scale): NSF CC*



Google Cloud Platform



Cisco Container Platform



Google Cloud



TACC Rodeo Cluster

TACC

TEXAS ADVANCED COMPUTING CENTER



A Lot Will Happen in Biology in the Next 15 Years

The tea leaves say...

In 2022, most funded biology labs are outsourcing DNA sequencing.

Terascale genomics experiments are common now and accessible to all in the US.

In 2032, every university research lab will have a DNA sequencer.

R1 research labs will move to the peta-/exascale in this PhD generation.

In 2037, all pharmacies, subways, hospitals, police stations, etc will have DNA sequencers.

These IoT DNA “sensors” will generate exabytes of data in aggregate each week.

We need more technical workers and start ups!

I am only talking DNA Sequencers...not CryoEM, Simulations, Medical Imaging



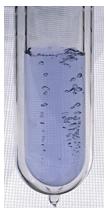
Basic intro to biological Information flow.

Organisms* are mostly 6 chemical elements (CHNOPS)

CARBON



HYDROGEN



NITROGEN



OXYGEN



PHOSPHOROUS



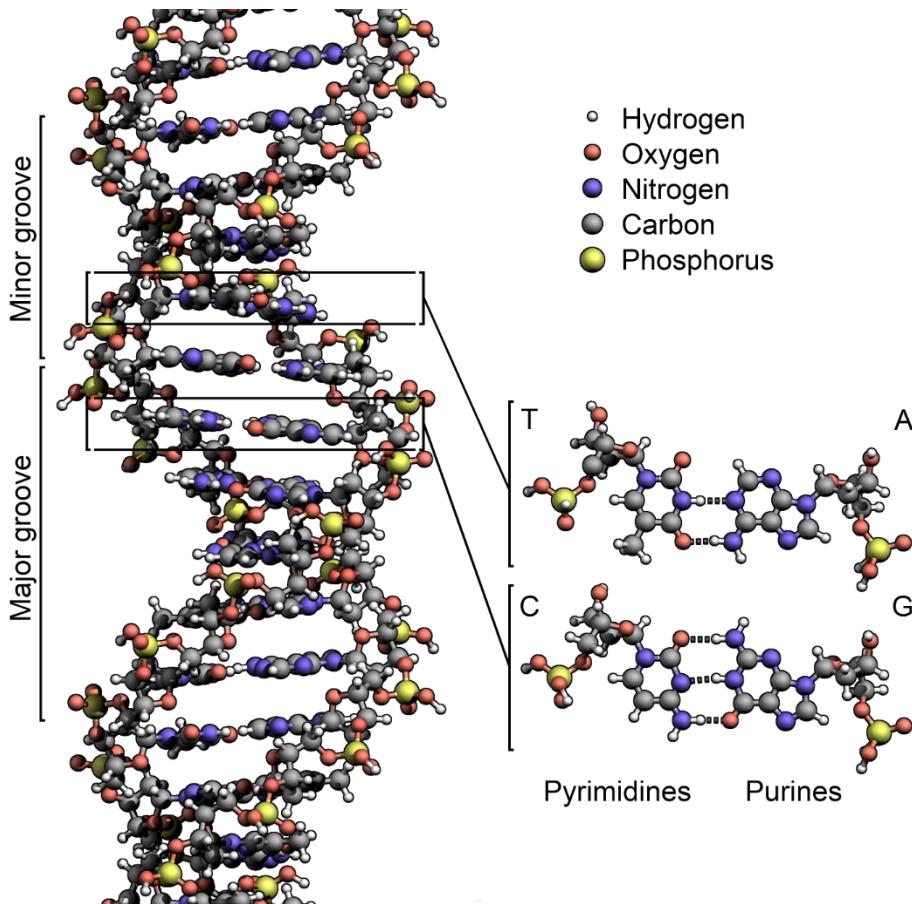
SULFUR



Group	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	
Period																			
1	H																	He	
2	Li	Be																	
3	Na	Mg																	
4	K	Ca	Sc	Ti	V	Cr	Mn	Fe	Co	Ni	Cu	Zn	Ga	Ge	As	Se	Br	Kr	
5	Rb	Sr	Y	Zr	Nb	Mo	Tc	Ru	Rh	Pd	Ag	Cd	In	Sn	Sb	Te	I	Xe	
6	Cs	Ba	*	Lu	Hf	Ta	W	Re	Os	Ir	Pt	Au	Hg	Tl	Pb	Bi	Po	At	Rn
7	Fr	Ra	**	103	104	105	106	107	108	109	110	111	112	113	114	115	116	117	118
*Lanthanoids		*	57	58	59	60	61	62	63	64	65	66	67	68	69	70			
**Actinoids		**	89	90	91	92	93	94	95	96	97	98	99	100	101	102			

*CHNOPS only applies to organisms from this region of the Milky Way.

DNA Is the Blueprint of Life and Built from CHNOP Elements



The alphabet of DNA are 4 molecules:

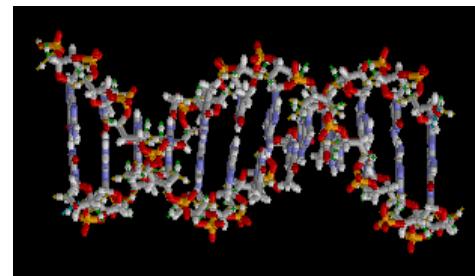
A
dene

T
yamine

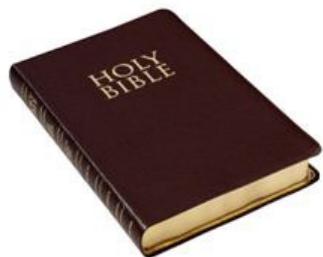
G
uanine

C
ytosine

The sequence of the alphabet makes words called genes that make the sentences, paragraphs, and chapters of an organism.



Genetics: Information in Human Book vs Human Genome

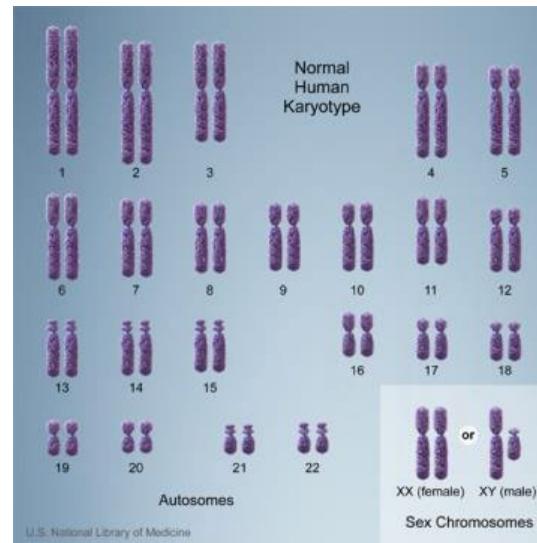


774,746 Words
26 letter alphabet

John 3:16: For God so loved the world that he gave his one and only Son, that whoever believes in him shall not perish but have eternal life.

John 3:16: For God so loved the world that he gave his one and only Son, that whoever believes in him shall not perish but have eternal life.

One letter change = big difference!



60,651 Words (Genes)
4 letter alphabet (ATGC)
6.4 billion letters

AATGGAGGCCACATAACACATTCAAACTTACTTGCAAAATAT

AATGGAGGCCACATAACACATGCAAACTTACTTGCAAAATAT

One letter change = higher risk for breast and ovarian cancer!

Change DNA and Change The Person: Coffee Example

chr15 (q24.1) 15p13 15p12 15p11.2 15q11.2 15q12(q13.1 q13.3 15q14 q15.1 15q21.1 15q21.2 15q21.3 15q22.2 15q22.3 15q23 q24.1 15q25.1 15q25.2 15q25.3 15q26.1 15q26.2 15q26.3



TGTGGGCACAGGAC



The **sequence of DNA contains information on how an organism responds to the world.**

TGTGGGCACAGGAC



TGTGGGC**CC**CAGGAC



TGTGGGCACAGGAC



TGTGGGCC**CC**CAGGAC



TGTGGGCC**CC**CAGGAC

Genotype	What It Means
  AA	Fast caffeine metabolizer: drinking coffee didn't increase subjects' heart attack risk
  AC	Slow caffeine metabolizer: drinking coffee increased subjects' heart attack risk.
  CC	Slow caffeine metabolizer: drinking coffee increased subjects' heart attack risk.

Information is Encoded as A,T,G,C; The Chromosome is The Storage Device

Blu-ray vs. DVD Capacity

Single-layer DVD

4.7 GB of data

average two-hour standard-definition movie
with a few extra features

Single-layer Blu-ray disc

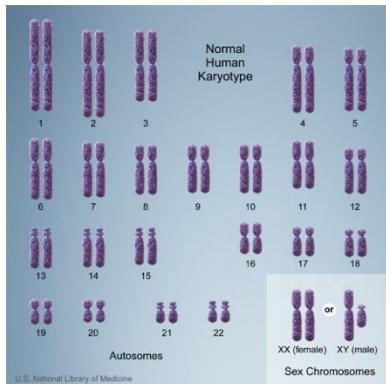
27 GB of data

more than 13 hours
of standard video

more than two hours of
high-definition video

01010111 01101001 01101011
01101001 01110000 01100101
01100100 01101001 01100001

Check out:
<http://www.ascii-code.com/>



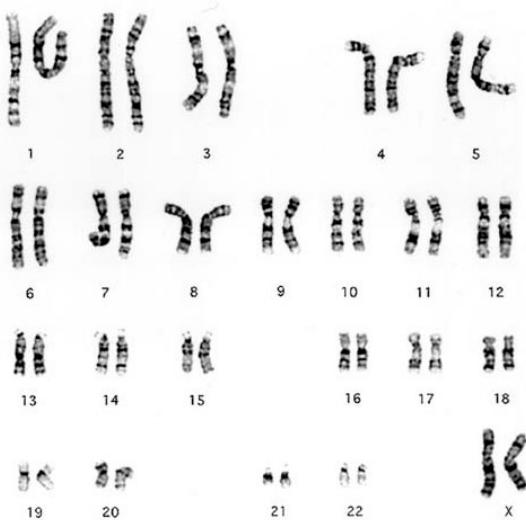
GGTTTCTATCTTACTAAAGGATATTCAGAAAATCTATATTCACTGAG
GAGGTGATAATTGGGTCTACTAACATTGAGACTGAAGTGGGCC
AGCAATAATTAAACTTATTATCCTTGAAGATTCAACTACTGGGGG
AGTACAACAGAAACAAATTAGAGAACATGAAGTTAATTACGTT
GAAGATGAAACATGGGACCCAACACTTGATCATTAG

All the information to make the unique version of us is encoded in
23 chromosome pairs spanning ~6,400,000,000 A,T,G,C base pairs (bytes)

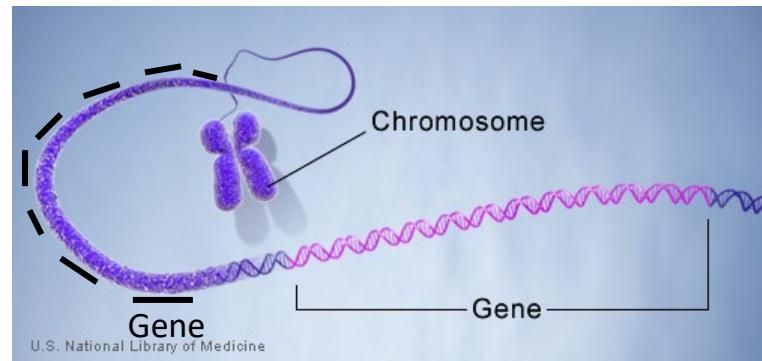
What is a genome? What is a gene?

Human Genome “Facts”

Initiation instructions for development
(3.2 billion*2) ATGCN
46 total chromosomes
24 unique chromosomes
(Chr1 = 250 million ATGCN)
60,651 genes



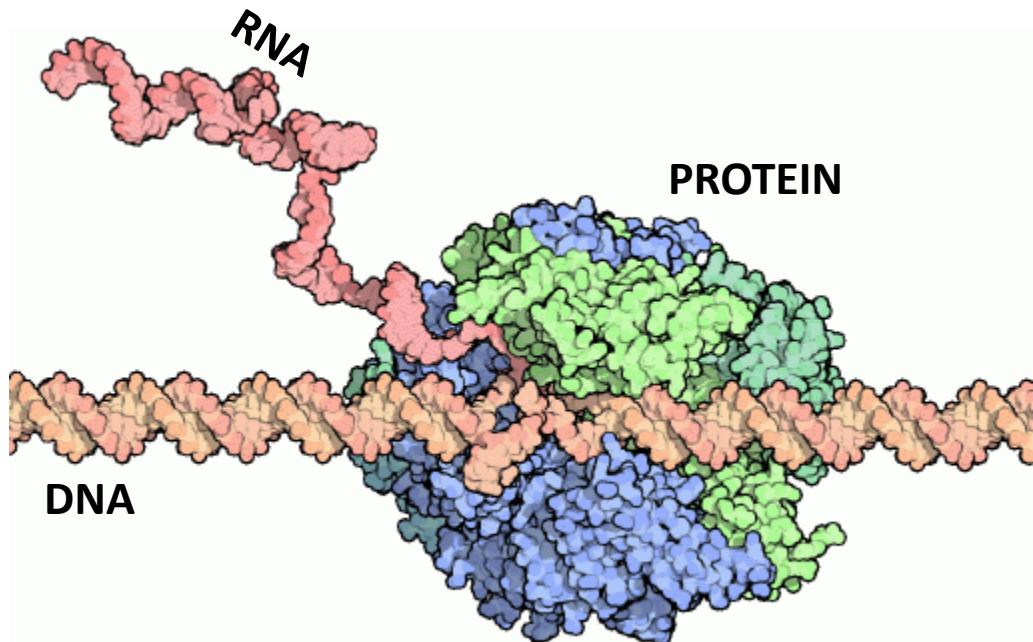
Courtesy of Dr. K. Phelan, Greenwood Genetic Center.
Noncommercial, educational use only.



- A **gene** is a **DNA interval** on the chromosome that gets copied (transcribe) into RNA.
- A **genome** is the set of **all chromosomes** and gene intervals.

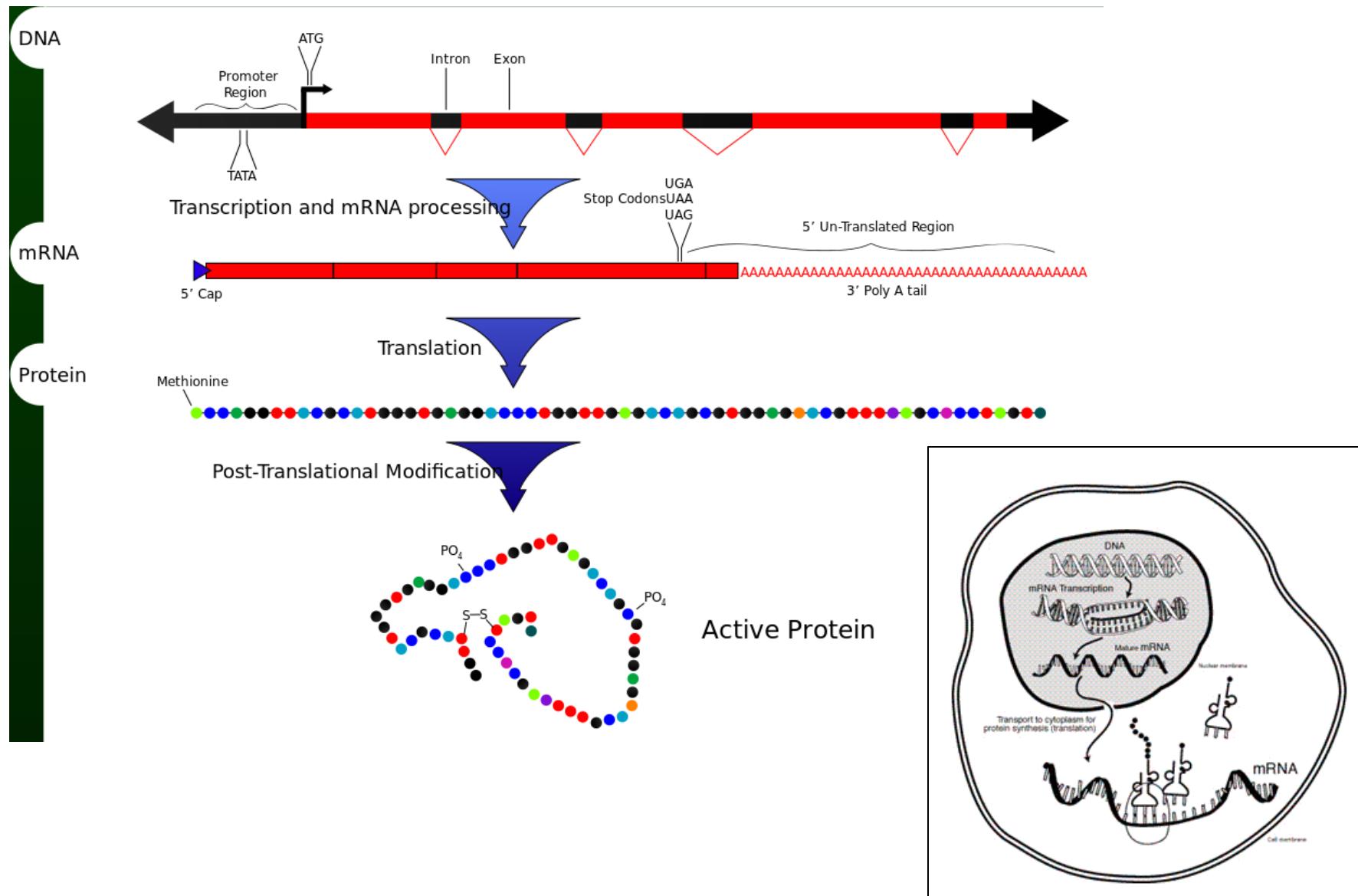


Chromosome Information is Accessed with Enzymes; Blu Ray Information is Accessed with Lasers



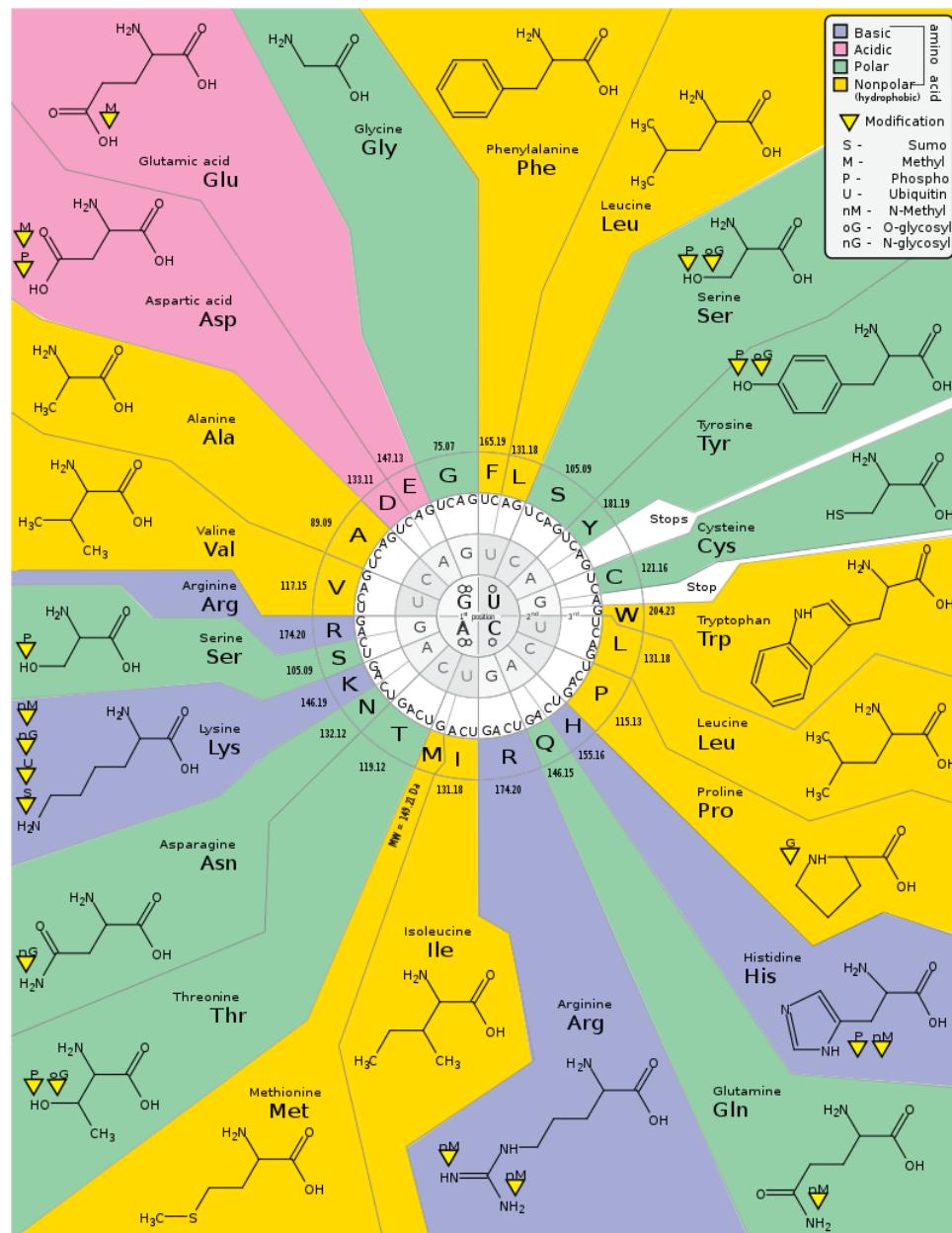
~The DNA molecule can be copied (Replication), transcribed (Transcription), and converted into protein (Translation)

Molecular Information Flow



Cell = Sacks in a sack

Molecular Information Flow Starts with the Genetic Code



>NBR1 Gene Fragment DNA (String of Nucleotides)
ATGGCAGTTAACAGGGAAACCAACTGCAGTGCAAG
TCCACGAAGGGCACCATGTCGTTGATGAAGCCCCAC



>NBR1 Gene Product RNA (String of Nucleotides)
AUGGCAGUUAAACAGGGAAACCAACUGCAGUGCAAG
UCCACGAAGGGCACCAUGUCGUUGAUGAAGCCCCAC



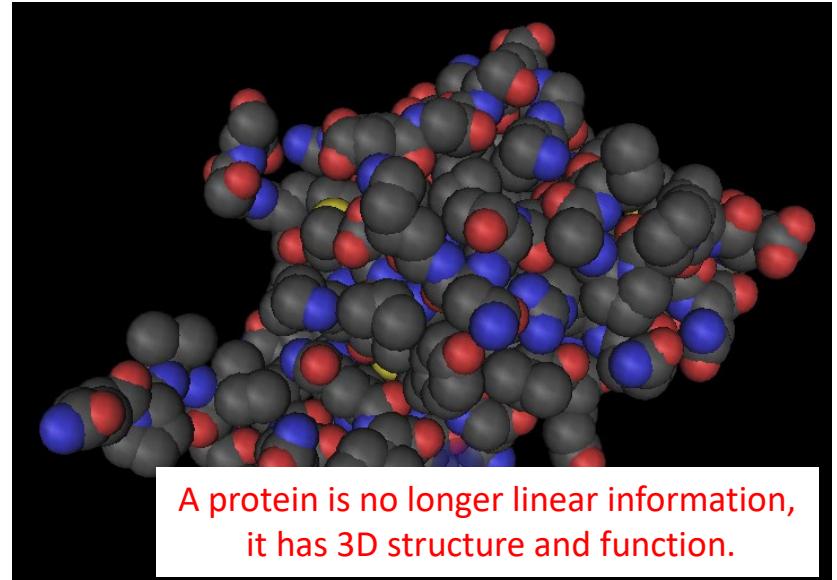
>NBR1 Gene Product Protein (String of Amino Acids)
MAVKQGNQLQCKSTKGTMSLMKPH

Second letter					
	U	C	A	G	
First letter	U	C	A	G	Third letter
U	UUU } Phe UUC } UUA } Leu UUG }	UCU } Ser UCC } UCA } UCG }	UAU } Tyr UAC } UAA Stop UAG Stop	UGU } Cys UGC } UGA Stop UGG Trp	U C A G
C	CUU } Leu CUC } CUA } CUG }	CCU } Pro CCC } CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } Arg CGC } CGA } CGG }	U C A G
A	AUU } Ile AUC } AUA } AUG Met	ACU } Thr ACC } ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G
G	GUU } Val GUC } GUA } GUG }	GCU } Ala GCC } GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } Gly GGC } GGA } GGG }	U C A G

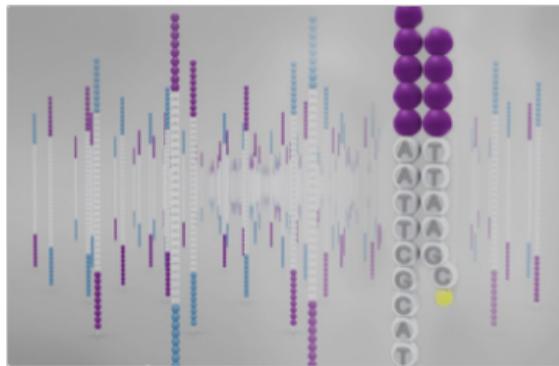
GGTTCTTACTAAAGGATATTCAAGAAAATCTATTCACTGAGGGAGGTGATAATTGGGCTACTAACATTGAGACT
 GAACGTAGGCCAGCAATAATTAAACTTATTACCTTGAGATTCACACTGGGGAGTACAACAGAAACAAATTAGA
 GAACATGAAGTTAATTACGTTGAAGATGAAACATGGGACCCACACTGATCATTTAGCTAACATGATGAGAAGAT
 GTACTTGGAAATAAGTGAACGAAAAGAAGATGGATTGAGATGGAGTAGAAGACAACAAATTGAAAGAGAATATGG
 AAAGAGCTGTTGATGCTTAGATATTACAGAACATGAACTCCAATTGGAACAGCAGTCAGGAAGAATATCTTAG
 TGATATTGCTATAAATCTAGGAGTACTAAATAAGAGGAAGCACATTAGTTAGTTAGTAGGTTCTGGCAGACTTATT
 CCCGTAAGAGACAGATAGTAAATATTAGGCTTGTGGACATACGTCCTGTCAGTAACTCAGTTCTGTTATTGTA
 GTGTGAAAGCAGCCACAGACAAATGTAACAAATGTGCTGTCTTCAATAAAACTTATCTACAAATACAGATAGTGGC
 CACATTGGCTGCGAGCTGAGTTGCTGACTTAGTGTAGATGACAGATTGTAGTGGAAAGATTTTAGTA
 TCAACAGACATTGAGGAGTTGCTATTATAATTCTGATGCTTTTCTATAATAGATGTCGTTCTGTTGGATTCTCT
 TCATATCTACTAGTGTACTCTCTGATTTGGGGTGTATTCTTACAGAATTCTTCTTCTTCTTCTTCTTCTT
 CTTTTTGAGAGAACAGGGTTGCTGTgATCACCCAGCTGGAGTGCAGTGGGGATCTGGCCCCCTGCGCTTCAACTC
 CTGGGCTTAAGCGATCTGTCACCTCAGCCTCCTAGTAGCTGGGACACAGGCTGCACCACTACACTGCTAATTG
 TATTTTTAGAGATGGGTTGCTGTGGAGACTGGTCTCACTGGCCCTGGCCAGAATTCTTACTTGTCTCATGTGAGAGAGTTCTAA
 TCCAAAGTGTGAAATTAGGTGAGTCAGTGTGCTTGTGACACTATAGTTTACTGTCTAGAGTCAGATTTTCTTAAATTCTTCA
 TAAACTCTGATGCTTGTGCTTGTGACACTATAGTTTACTGTCTAGAGTCAGATTTTCTTAAATTCTTCA
 TTTCTTAATTCTCATTGTTGCTTAGTTACAAATTACTTAAGATGCCTATATTCTGAAATTATATTCTTCA
 ATAAATAAGTTAAAAAATTAGATTGACCATTAAATTAGAGTTATTCCCCATCTTCTAATGGCAACATAATTG
 ATTGAGCAATTCTCTGATTGGTAACTATTAGTAATCACTTTATCAGTATAAACATGCAATTCTTAAATTGGTTTC
 AAAATATTCTCATACCAACCATTCTAA
 ATTTCTAGGAATTCTTATTAGGTATAC
 TTTTTTGTGAGATGGAACTATCTCTG
 CCGGTTCAAGTGAATTCTCACCTCAGG
 ATTTTAGTAGAGAAAGTGTGGGATT
 TACACCAATTATTGCCATTTTTATTGG
 TTAAAAAATAAAACTTTTCTAGTACAAAGCTTATTCAATTTTTATTGCTTAAATTGGCAGGTTAAATTAGGC
 TGACTCTTTCTACTCTGTTAAATTGGCATTTTGAGGGTTCTTCTATGGATTGGGGATCCCAC
 CGTTGGACACTGGCTCTGGAGATTG
 TAACATCTGACATTATTCAATTATAGAGAATGCTTGTAAAATACTGATTATTAAAGATATTACATCATAG
 AGTTTTGCAATCTTAAATCAGTGGATG
 GTCAGGAACAACATTTGTTGACGGCTTCAACATTGTGAAG
 GTTCTTACTCAAAAGTTGGAGACACAGTGTGACATTGTTAACCTGAGTTGACATTGTCTTGTATTCTGTTAG
 AGCAGGGATTATTGCTCTTATCAACTTGTACTTCAACTTCTTCTTCTTGGGTTCTTCTTGTCT
 TCTTCTCTCTCCCTTTTATCATATAATATGAAATATATTATATCATTCATTAGGGAGAGGAAGTTGCTAA
 GATATCTAGTATAGGAGCTTGTCTGAGAAGGAGATGAGTC
 TTAAGCATTATCTCCAATGATAATG
 CTAAGGTATGTTACAATTATAAAATTACTCAAGTCTTCAAAGGACATTAAATTAGTAAATTAAACTAATTCT
 AACTAGGTTCTACCAATGAAATTGCTACTAATTGTAACATTAGATTCACTTCAATTCTCATGTTCTTCTCATGTAG
 TCTATAAAATTGGGTTAGGAGTAATTACTAATTAAATTGCTTCTTGTCTTCTTATTGAAAGACAGGG
 TCTCACTCTACCCAGGCTGGAGTGCAGTTGCTCCATCTCGGCTACTGCAACCTCACCTCTGGCTCAAGTGTACTC
 CTGCATAAGCCTCCGAGTAGCTGGGATTGGCGTCACCCATGCCGGCTAATTGTTAGTTAGAGACAG
 GGTCTCACCATGTTGCTTGTCTGCAACTCTCAAGTAATCCACCCGGCTTGGCTCCAAAGTGTGGGATT
 ACAGGTGAGGCCACACCTGGCTGGATTCTGAGTTGACTGACATTTACACAAATTATTGCCATT
 GGTACCAATTAGTATTGTTCTTATTGCTTAAATAAGGATAAAACTTTTTAAAGAACACTTTCTTAG
 TAAAAGCTTATTGCTTAAATTGCTTAAATTGAGGTTCTTCTGAGGTTAAAGGCTGCTCTTCTCTTCTGTTAA
 TGGTGGATTGGGTTCTTCTATGGATTGGGGATCCCATCTGGGACACTTGGTCTGGAGATTG
 CCAAATTATTGAAAGTTCTCATACCAACATTGGCAGTAGACTGAAATAAAACATCTGACACTATTCTTCAATT
 ATAGAGAAATGCTATTGTAATTACTGATTATTAAAGATATTACAGTCAGTAGTTTGTCAATCTTAAATCAGTGG
 ATGTCAAAACAGCTATTGTTCTGACGATTTCACATTGTAAGAGTCTTACTACAAAGTTGGAGACACA
 GTTGATATTGTTAACCTGAGTGAACCTGTCATTGTTATTCTGTTAGAGCAGGGATTATTGTTCTTCT
 TTGTTCTTCTTCTTCTTGTCTTGTCTTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCTCT
 GTAAATATATTATATCCATTAGGGAAGAGGAAGCTGCTAAAGATATTAGTATAGGAGCTTGTGAGAAA
 GAGGATATGAAGTCAATTATGGAAATTAAATGCTTAACATTGTTAAAGCATTATCTCCAATGATAATGAAACGA
 TAGTCCTATGTAATTGAGAGTGTAGAAGATTAGAAATGGAGATGCTTAAG

GUUUCUAUCUUACUAAGGAUAAAUCAGAAAUCUAUAU
 CACUGAGGAGGAUGAUAAAUGGGCUACUAACAUUGAGAC
 GAACUGAGGCCAGCAUAUUUAACUUUAUCCUUUGA
 AGAUUCAACUACUGGGGAGUACAACAGAAACAAU
 AACAGAAGUUUL
GENE PRODUCT
 JGGAGAAGAUGU
 RNA="AUGC"
 (information)
 AUGGAUUGAAG
 AAUAUGGAAAGA
 AACAUAGAACUCCA
 AAUUUUGGAACAGCAGUCUCAGGAAGAAUAUCUUAGUGAUA
 UUGCUUAAAUCUACUGAGCAUUUAUCUCCAAUGAU
 GAAAACGAUACGUCCUAUGUAAUUGAGAGUGAUGAAGAU
 AGAAAUGGAGAUGC
 UUAAG

MIIGSTNIETLPSNNLNL
 SFE
 DTTGGVQQKQIRHEV
 LIHVEDETWDP
 TLDHLAKHDGEDVLGNK
 VERKEDGF
 FEDGV
 EDNKL
 KENMERAC
 LMSLDI
 EM
GENE PRODUCT
 LK
 PROTEIN = "20 unique amino acids letters"



How does one obtain DNA sequences from a sample?

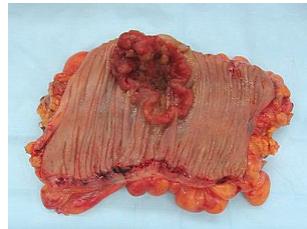


How Does Illumina NGS Work?

Illumina sequencing utilizes a fundamentally different approach from the classic Sanger chain-termination method. It leverages sequencing by synthesis (SBS) technology – tracking the addition of labeled nucleotides as the DNA chain is copied – in a massively parallel fashion.

Next-generation sequencing generates masses of DNA sequencing data, and is both less expensive and less time-consuming than traditional Sanger sequencing.² Illumina sequencing systems can deliver data output ranging from 300 kilobases up to multiple terabases in a single run, depending on instrument type and configuration.

Sample



Purify DNA



DNA Sequencer



Computational Analysis



DNA SEQUENCE/QUALITY FASTQ FILE

(10 – 1000 Million “DNA Reads”; \$slen=100-300 characters)

```

@SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
GGGTGATGCCGCTGCCGATGGCGTCAAATCCCACC ←
+SRR001666.1 071112_SLXA-EAS1_s_7:5:1:817:345 length=36
IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII9IG9IC ←

```

SEQUENCE
QUALITY

http://en.wikipedia.org/wiki/FASTQ_format

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	Ø	96	60	140	`	`
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A	97	61	141	a	a
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B	98	62	142	b	b
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C	99	63	143	c	c
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D	100	64	144	d	d
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E	101	65	145	e	e
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F	102	66	146	f	f
7	7	007	BEL (bell)	39	27	047	'	!	71	47	107	G	G	103	67	147	g	g
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H	104	68	150	h	h
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I	105	69	151	i	i
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J	106	6A	152	j	j
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K	107	6B	153	k	k
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L	108	6C	154	l	l
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M	109	6D	155	m	m
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N	110	6E	156	n	n
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O	111	6F	157	o	o
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P	112	70	160	p	p
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q	113	71	161	q	q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R	114	72	162	r	r
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S	115	73	163	s	s
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T	116	74	164	t	t
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U	117	75	165	u	u
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V	118	76	166	v	v
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W	119	77	167	w	w
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X	120	78	170	x	x
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y	121	79	171	y	y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z	122	7A	172	z	z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[123	7B	173	{	{
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\	124	7C	174	|	
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]]	125	7D	175	}	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^	126	7E	176	~	~
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_	127	7F	177		DEL

Source: www.LookupTables.com

http://en.wikipedia.org/wiki/Phred_quality_score
<http://www.asciiitable.com/>

Base Call Quality (Q) score is
ASCII Decimal# minus 33
(usually).
 $Q(I)=40:::99.999\%$ Accuracy

$$Q = -10 \log_{10} P$$

or

$$P = 10^{-\frac{Q}{10}}$$

Phred quality scores are logarithmically linked to error probabilities

Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%
50	1 in 100000	99.999%

What do we do with DNA sequences?

(Align \$little.string to \$big.string)

Sequencing Output Per Flow Cell*

NovaSeq 6000 System

Flow Cell Type	SP	S1	S2	S4
1 × 35 bp	N/A	N/A	N/A	280-350 Gb
2 × 50 bp	65-80 Gb	134-167 Gb	333-417 Gb	N/A
2 × 100 bp	134-167 Gb	266-333 Gb	667-833 Gb	1600-2000 Gb
2 × 150 bp	200-250 Gb	400-500 Gb	1000-1250 Gb	2400-3000 Gb
2 × 250 bp	325-400 Gb	N/A	N/A	N/A

*Specifications based on Illumina PhiX control library at supported cluster densities.

```
for (my $i=0; $i <= $dna.read.count; $i++) {  
    for (my $k=0; $k <= $chromosome.count; $k++) {  
        return_genome_coordinate ($DNA[$i], $chromosome[$k])  
        if ($coordinate overlaps $gene) {$gene.$count++}  
        else {evidence.for.new.gene}  
    }  
}
```

chr	total length
1	249,250,621
2	243,199,373
3	198,022,430
4	191,154,276
5	180,915,260
6	171,115,067
7	159,138,663
8	146,364,022
9	141,213,431
10	135,534,747
11	135,006,516
12	133,851,895
13	115,169,878
14	107,349,540
15	102,531,392
16	90,354,753
17	81,195,210
18	78,077,248
19	59,128,983
20	63,025,520
21	48,129,895
22	51,304,566
X	155,270,560
Y	59,373,566

Example Alignment of 1 out of 2×10^8 Sequences

Primary Assembly	
chr	total length
1	249,250,621
2	243,199,373
3	198,022,430
4	191,154,276
5	180,915,260
6	171,115,067
7	159,138,663
8	146,364,022
9	141,213,431
10	135,534,747
11	135,006,516
12	133,851,895
13	115,169,878
14	107,349,540
15	102,531,392
16	90,354,753
17	81,195,210
18	78,077,248
19	59,128,983
20	63,025,520
21	48,129,895
22	51,304,566
X	155,270,560
Y	59,373,566

Length of \$chromosome.01

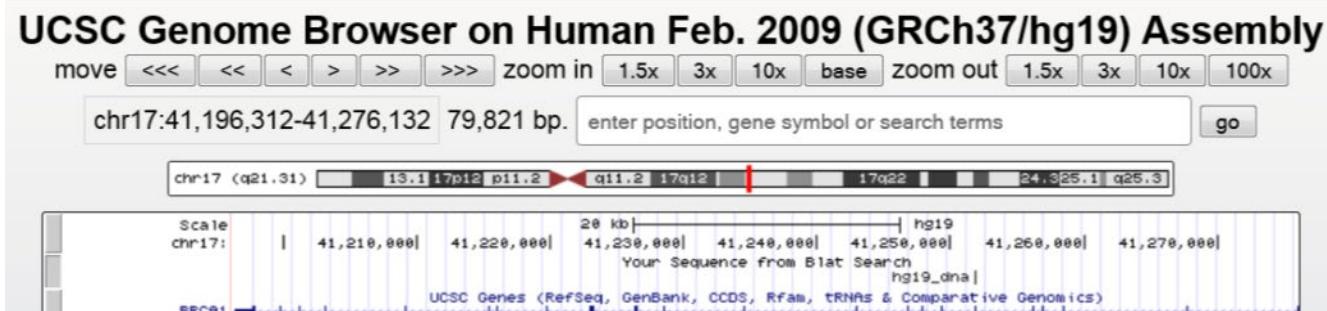
Read01 out of 200,000,000:

TCAATGTAGA CAGACGTCTT TTGAGGTTGT ATCCGCTGC

Search ALL chromosomes for alignment to Read01:
Hit! on Chromosome 17 @ 41,251,797-41,251,835

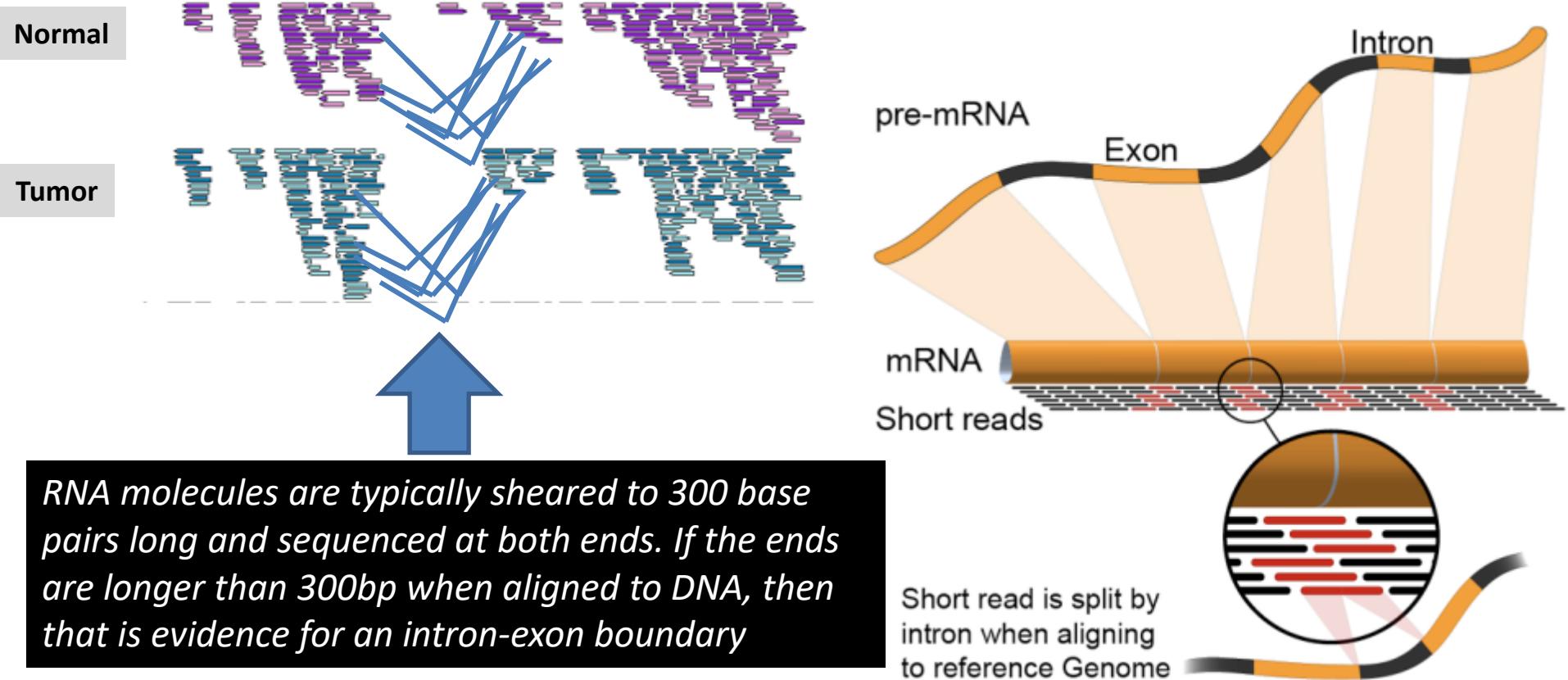
00000001 tcaatgtagacagacgtctttgagggttatccgctgc 00000039
>>>>> ||||||| ||||||| ||||||| ||||||| >>>>>
41251797 tcaatgtagacagacgtctttgagggttatccgctgc 41251835

Add one to number of BRCA1 RNA transcripts expressed
Since this gene is encoded in the DNA at that coordinate.



Application: Counting RNA Molecules for Detection of Differential Gene Expression (DEG)

Are molecule counts significantly different at any given gene between experimental conditions?



<http://www.data2bio.com/services/rnaseq/>

<http://en.wikipedia.org/wiki/RNA-Seq>

OUTPUT: Sample X Gene (Gene Expression Matrix: GEM)

Genes

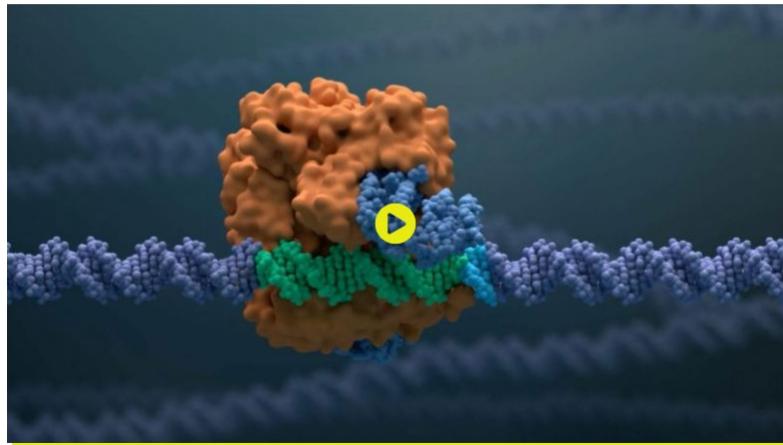
Samples

	ATLAS	C.AMBER	GRASSL	LEOTI	PI_267573	PI_506030	PI_506069	PI_510757	RIO	WRAY	PI_229841	PI_297130	PI_329311
Sobic.001G000100.1	43	80	59	57	63	55	141	81	98	64	39	23	34
Sobic.001G000200.1	5019	4612	4159	4997	5025	5037	5688	5192	4721	4790	5003	5318	4080
Sobic.001G000300.1	13	211	363	197	816	16	22	25	609	33	936	1008	812
Sobic.001G000400.1	96	247	155	225	173	181	176	122	21	88	148	142	111
Sobic.001G000400.2	18214	16957	20109	22401	22026	14420	15878	17051	20284	21036	22790	24456	21542
Sobic.001G000500.1	0	2	0	0	0	0	0	2	0	0	0	0	0
Sobic.001G000600.1	158	895	1053	70	1405	1247	1264	1036	1194	35	70	176	1931
Sobic.001G000700.1	362	175	533	145	924	133	199	224	436	256	343	301	598
Sobic.001G000700.2	9223	7622	7228	8348	8610	10394	9611	10517	10360	8733	11743	11149	13108
Sobic.001G000800.1	13066	13519	12412	13375	14575	14003	15294	12427	14094	14013	13180	13968	12237
Sobic.001G000900.1	35	50	44	28	27	32	21	12	22	16	50	18	11
Sobic.001G001000.1	268	279	100	222	256	424	490	202	222	241	131	103	14
Sobic.001G001100.1	32	38	14	28	24	89	89	52	41	55	44	29	8
Sobic.001G001200.1	128	146	183	165	144	159	284	223	178	190	282	248	343
Sobic.001G001300.1	709	670	1114	1124	971	889	1092	826	809	909	926	395	753
Sobic.001G001400.1	12837	11848	10055	12620	13128	13101	12425	11465	12352	12017	12296	13688	8615
Sobic.001G001500.1	7002	6330	5217	6856	6767	6375	7224	5941	6374	6656	6639	8265	5347
Sobic.001G001600.1	0	0	0	0	0	17	0	9	0	0	0	29	0
Sobic.001G001700.1	83	0	11	49	43	76	0	89	0	103	95	211	0
Sobic.001G001700.2	33	1	17	11	0	0	0	1	8	13	27	0	0
Sobic.001G001800.1	3	0	0	3	6	0	0	6	0	0	10	29	0
Sobic.001G001900.1	5550	5800	4852	4827	4740	5997	4939	4807	4943	4973	5813	5272	5730
Sobic.001G002000.1	9862	11281	9520	10405	12690	11908	12670	10716	10612	10418	11096	12666	9160
Sobic.001G002100.1	1253	984	1615	1517	1244	1064	973	1064	1626	1675	746	1237	914
Sobic.001G002200.1	0	0	0	0	0	0	0	0	0	0	0	0	0
Sobic.001G002300.1	31	44	58	26	24	53	24	13	24	25	10	31	0
Sobic.001G002400.1	77	24	66	11	15	44	26	12	115	44	2	99	0
Sobic.001G002500.1	38	17	11	48	57	6	42	41	98	22	16	37	11
Sobic.001G002600.1	2189	1300	3272	3372	3617	4483	3582	3795	3332	2515	3798	3413	2815

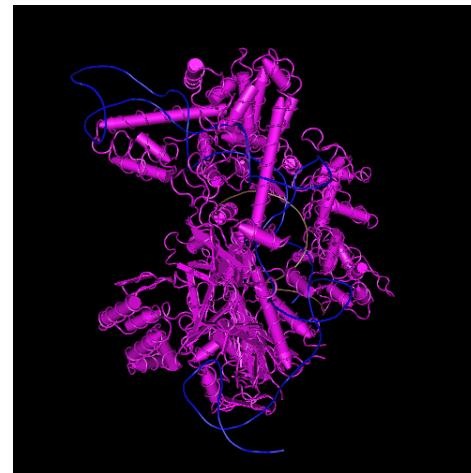
The values are molecules counted at each gene interval.

Physically editing DNA sequences with CRISPR/CAS9

CRISPR/Cas9



STRUCTURE:



Current Opinion in Structural Biology
Volume 69, August 2021, Pages 86-98



Engineering Cas9 for human genome editing

Ian M. Slaymaker, Nicole M. Gaudelli

Show more ▾

+ Add to Mendeley Share Cite

<https://doi.org/10.1016/j.sbi.2021.03.004>

[Get rights and content](#)

Questions:

How do you know which genes to edit?

How do you get the CAS9-RNA complex to cells?

OUTPUT: SAM/BAM Alignment File

1.1 An example

Suppose we have the following alignment with bases in lower cases clipped from the alignment. Read r001/1 and r001/2 constitute a read pair; r003 is a chimeric read; r004 represents a split alignment.

```
Coor      12345678901234  5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT

+r001/1    TTAGATAAAGGATA*CTG
+r002    aaaAGATAAA*GGATA
+r003    gcctaAGCTAA
+r004        ATAGCT.....TCAGC
-r003        ttagctTAGGC
-r001/2        CAGCGGCAT
```

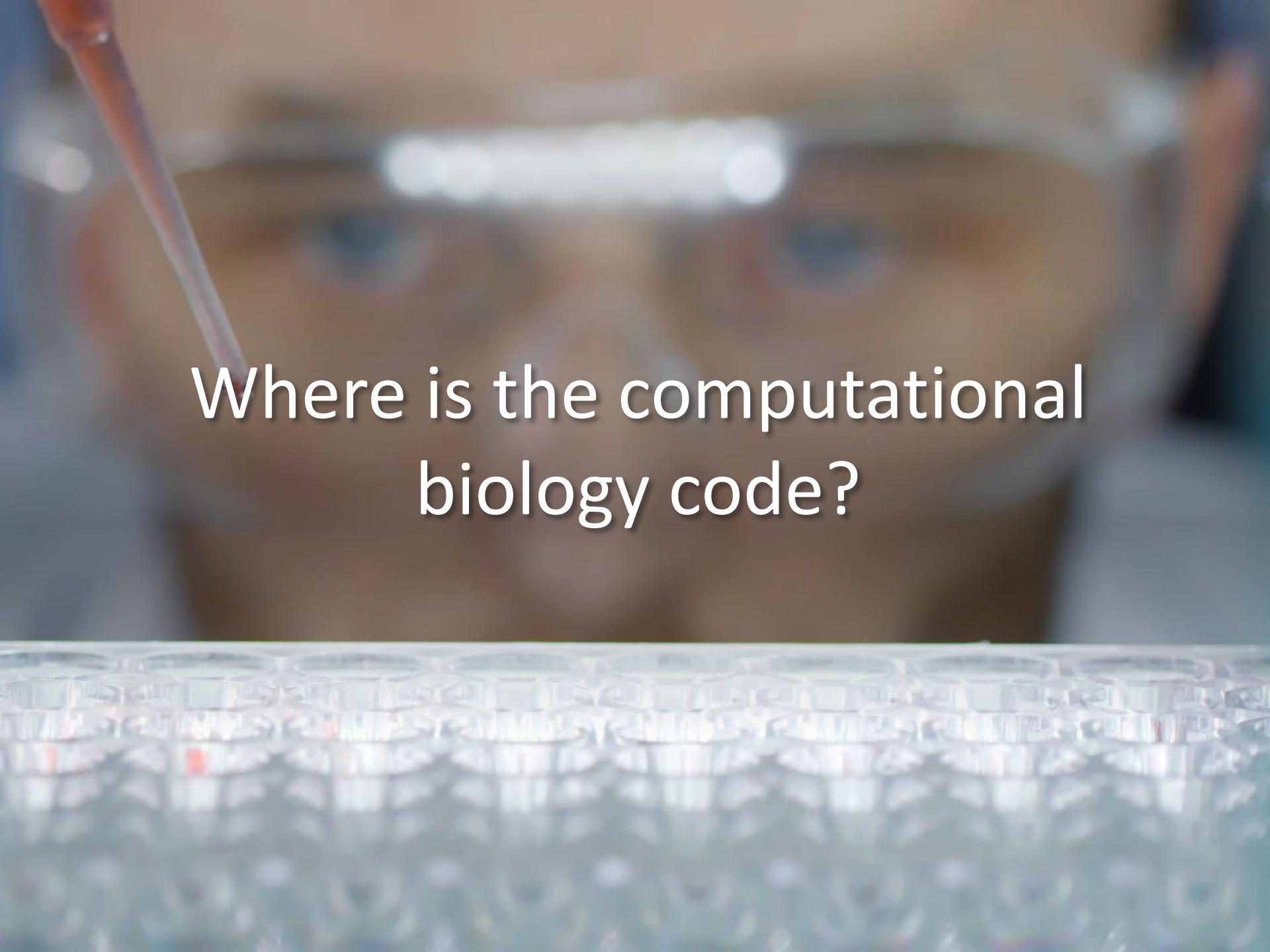
MILLIONS OF ALIGNMENTS!

SAM = ASCII
BAM = Binary

The corresponding SAM format is:

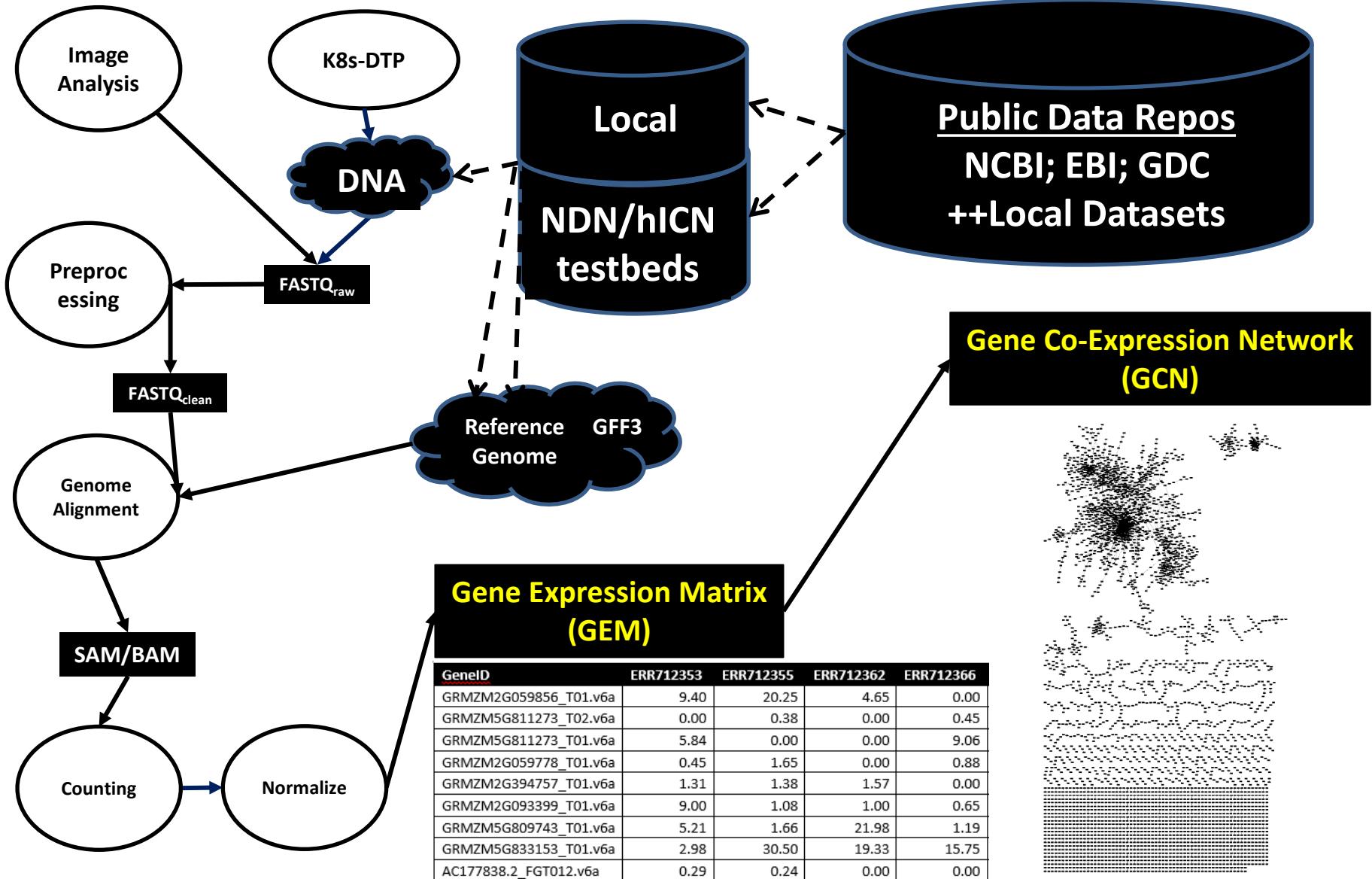
```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001  99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002  0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003  0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004  0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

SPECIFICATION::: <http://samtools.github.io/hts-specs/SAMv1.pdf>



Where is the computational
biology code?

Computational Workflows Are Part of the Modern Biology Lab



nextflow

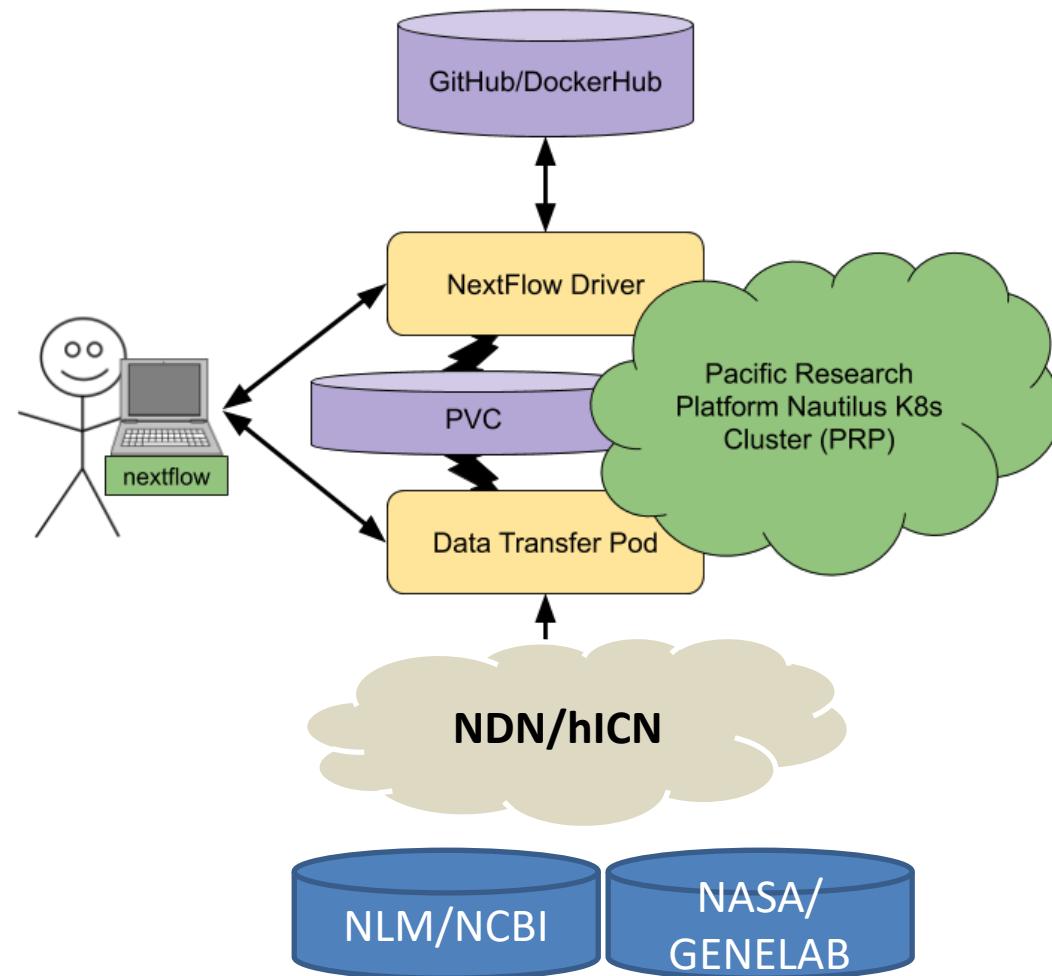
Galaxy



GEMmaker

We Wrap Our Containerized Genomic Workflows in NextFlow

Nextflow is a workflow manager that enables the streamlined execution of multiple steps in a scientific workflow. Users can simply change a configuration file to modify a workflow.



Workflows:

- [SystemsGenetics/KINC-nf](#)
- [SystemsGenetics/GEMmaker](#)
- [SystemsGenetics/GEMprep](#)
- [SystemsGenetics/gene-oracle](#)

Additional Tools:

- [SystemsGenetics/kube-runner](#)
- [SciDAS/nextflow-api](#)

Kubernetes Platforms:

- Pacific Research Platform (PRP)
- Cisco Container Platform (CCP)
- TACC Rodeo
- Google Cloud Platform (GCP)

Data-intensive compute (DIC) workflow development is impossible in the commercial cloud for small labs.
Democratized DIC systems are CRITICAL.

GEM Construction with GEMmaker Containerized Workflow

DOI [10.5281/zenodo.1283750](https://doi.org/10.5281/zenodo.1283750)

build passing

<https://github.com/SystemsGenetics/GEMmaker>



GEMmaker is a [Nextflow](#) workflow for large-scale gene expression sample processing, expression-level quantification and Gene Expression Matrix (GEM) construction. Results from GEMmaker are useful for differential gene expression (DGE) and gene co-expression network (GCN) analyses. The GEMmaker workflow currently supports Illumina RNA-seq datasets.

GEMmaker uses the following tools:

- [python3 v3.5.1](#)
- [nextflow v0.32](#)
- [sratoolkit v2.9.2](#)
- [fastQC v0.11.7](#)
- [trimmmomatic v0.38](#)
- [hisat2 v2.1.0](#)
- [samtools v1.3.1](#)
- [stringTie v1.3.4d](#)
- [MultiQC v1.5](#)



GEMmaker is a collaborative project of the [Ficklin](#) and [Feltus](#) programs at Washington State University and Clemson University respectively with guidance from [RENCI](#).

GEMmaker is funded by the [NSF SciDAS](#) project, award #1659300



Where is the data?

Public Data Source: NCI Genome Data Commons

11,315 Samples; 33 Tumor Sub-Types + “Solid Tissue Normal”; 10,610 Pubmed Hits

The screenshot shows the homepage of the GDC Data Portal. At the top, there's a navigation bar with tabs for Home, Projects, Exploration, Analysis, and Repository. Below the navigation is a search bar with placeholder text "e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2". To the right of the search bar is a large graphic featuring a human silhouette with internal organs highlighted in various colors (purple, yellow, green). To the right of the silhouette is a bar chart titled "Cases by Major Primary Site". The chart lists 33 tumor sub-types and their corresponding case counts. A callout box highlights the Lung category, which has 12,026 cases (64,330 files). Below the main header, there's a "Data Portal Summary" section with counts for Projects (67), Primary Sites (68), Cases (84,392), Files (596,758), Genes (23,399), and Mutations (3,287,299). At the bottom of the page, there's a section titled "GDC Applications" with links to various tools like Data Portal, Website, API, Data Transfer Tool, Documentation, Data Submission Portal, Legacy Archive, and Publications. The footer contains links to Site Home, Policies, Accessibility, FOIA, Support, and various US government agencies. It also includes a note about NIH's mission and the release date.

tcga - Google Search X GDC

portal.gdc.cancer.gov

National Cancer Institute GDC Data Portal

Home Projects Exploration Analysis Repository

Quick Search Manage Sets Login Cart 0 GDC Apps

Harmonized Cancer Datasets

Genomic Data Commons Data Portal

Get Started by Exploring:

Projects Exploration Analysis Repository

e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

Data Portal Summary Data Release 27.0 - October 29, 2020

PROJECTS 67 PRIMARY SITES 68 CASES 84,392 FILES 596,758 GENES 23,399 MUTATIONS 3,287,299

Cases by Major Primary Site

Tumor Sub-Type	Number of Cases
Adrenal Gland	1
Bile Duct	1
Bladder	1
Bone	1
Bone Marrow	9
Brain	1
Breast	9
Cervix	1
Colorectal	8
Esophagus	1
Eye	1
Head and Neck	1
Kidney	3
Liver	2
Lung	12,026 cases (64,330 files)
Lymph Nodes	1
Nervous System	1
Ovary	3
Pancreas	2
Pleura	1
Prostate	1
Skin	2
Soft Tissue	1
Stomach	1
Tie	1
Thymus	1
Thyroid	1
Uterus	2

GDC Applications

The GDC Data Portal is a robust data-driven platform that allows cancer researchers and bioinformaticians to search and download cancer data for analysis. The GDC applications include:

Data Portal Website API Data Transfer Tool Documentation Data Submission Portal Legacy Archive Publications

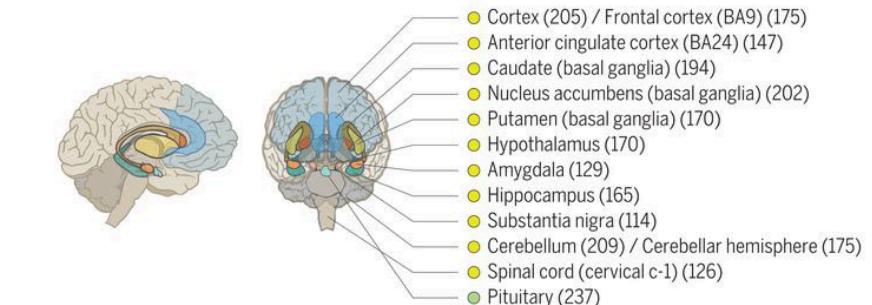
Site Home | Policies | Accessibility | FOIA | Support
U.S. Department of Health and Human Services | National Institutes of Health | National Cancer Institute | USA.gov
NIH... Turning Discovery Into Health ©
UI v1.27.0 @ 158fc606, API v3.0.0 @ dfa39447, Data Release 27.0 - October 29, 2020

Type here to search

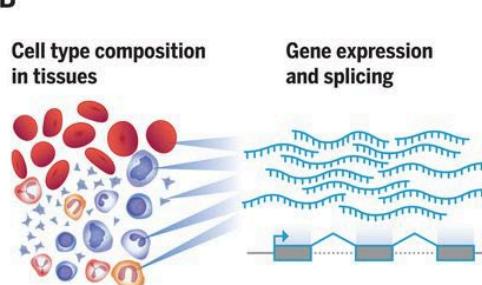
Source: <https://portal.gdc.cancer.gov/>

Public Data Source: Genotype-Tissue Expression (GTEx) Project

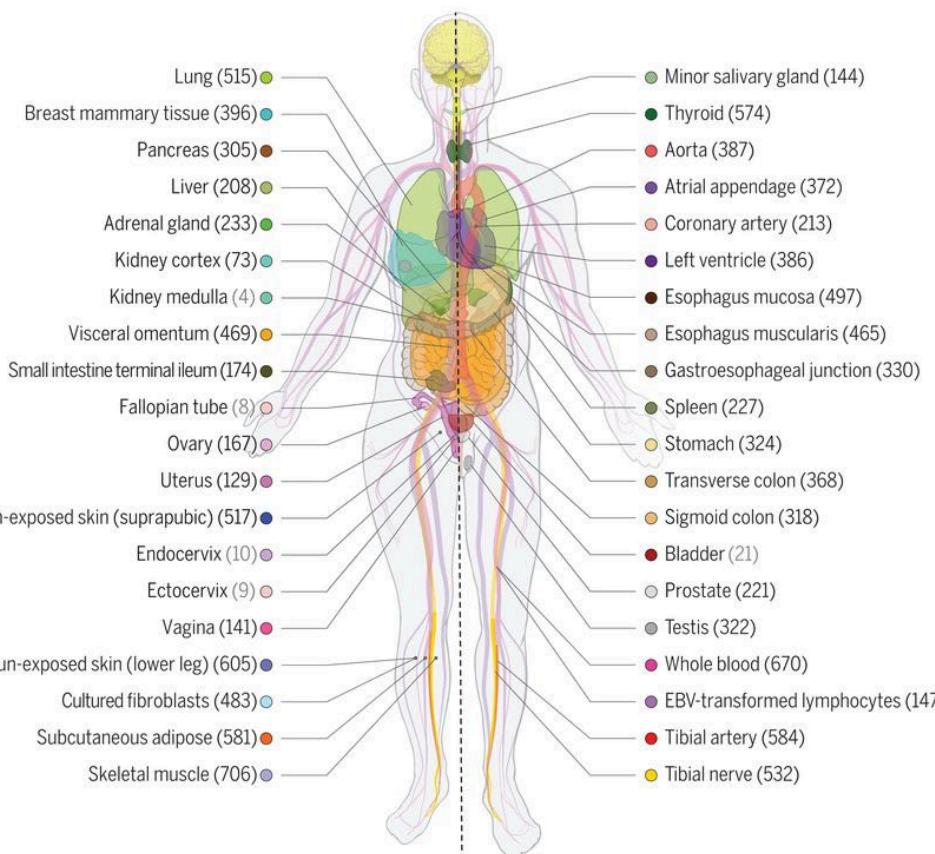
A



B

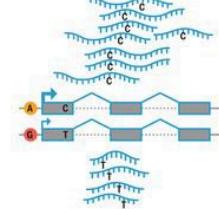


Release V8
838 donors
49 tissues
15,201 RNAseq



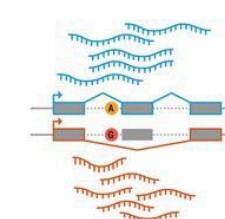
Expression quantitative trait loci (eQTLs)

cis-eQTLs

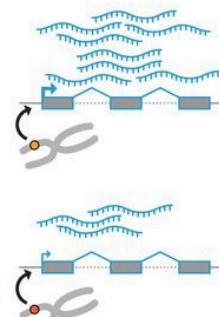


Splicing quantitative trait loci (sQTLs)

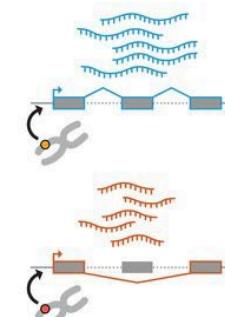
cis-sQTLs



trans-eQTLs



trans-sQTLs



The GTEx Consortium atlas of genetic regulatory effects across human tissues

The GTEx Consortium*

* See all authors and affiliations

Science 11 Sep 2020;
Vol. 369, Issue 6509, pp. 1318-1330
DOI: 10.1126/science.aaz1776

A photograph showing several business people in suits gathered around a table, looking at a large tablet device. One person's hands are visible on the tablet screen, which displays a data visualization with a large orange circle. There are also smartphones and coffee cups on the table.

What do you do
with the data?

Example Publication

Patterns

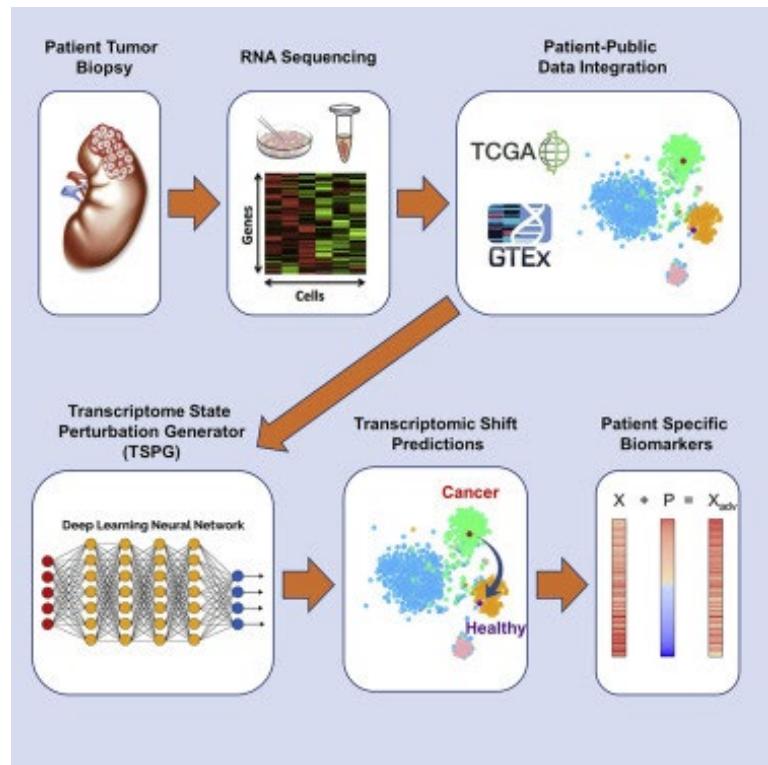
Volume 1, Issue 6, 11 September 2020, 100087



Article

Cellular State Transformations Using Deep Learning for Precision Medicine Applications

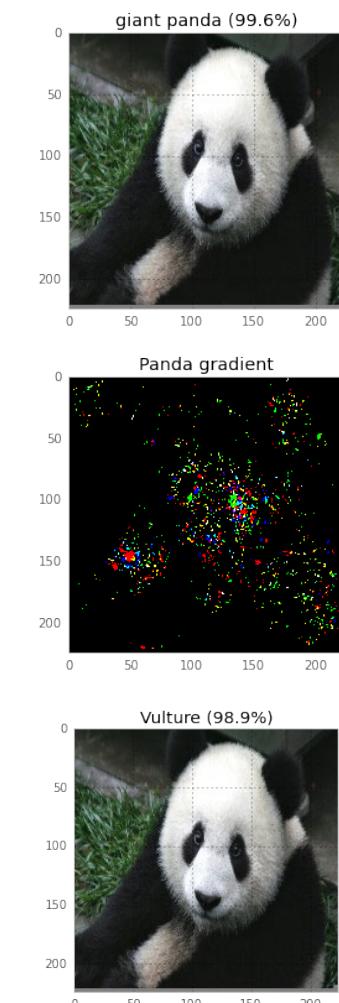
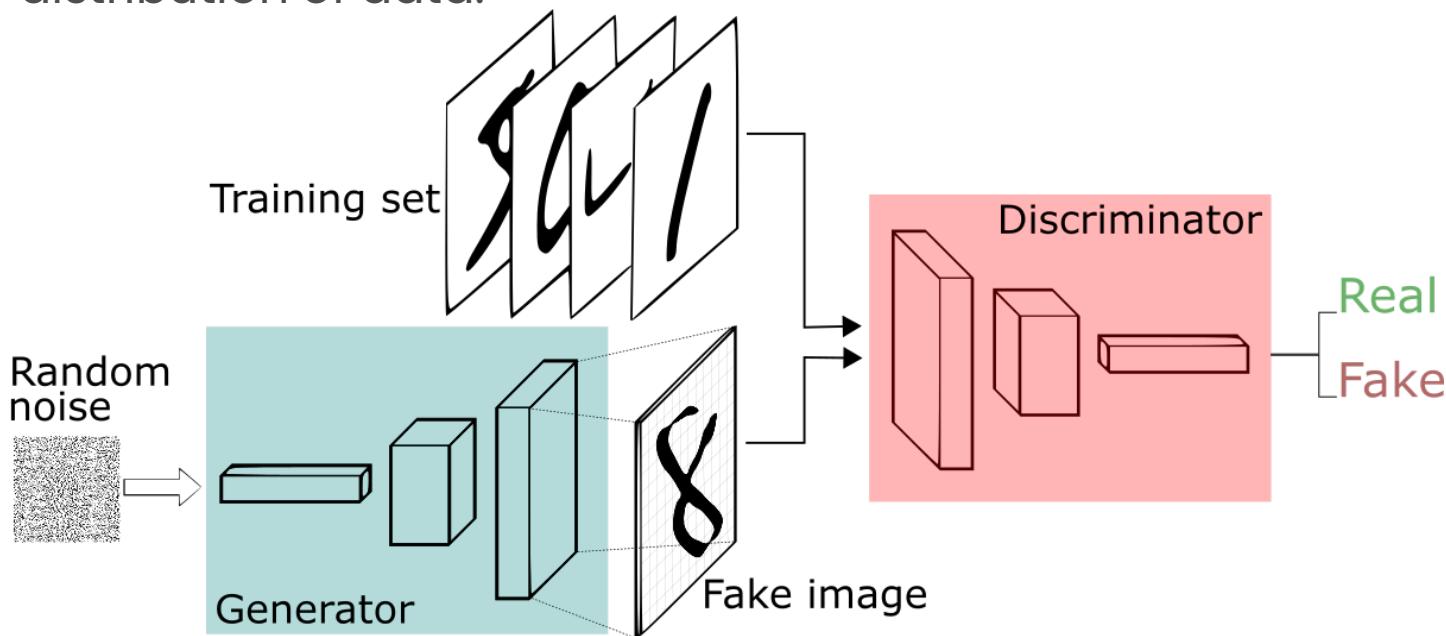
Colin Targonski ^{1, 6}, M. Reed Bender ^{2, 6}, Benjamin T. Shealy ¹, Benafsh Husain ², Bill Paseman ³, Melissa C. Smith ¹, F. Alex Feltus ^{2, 4, 5, 7}



<https://github.com/ctargon/TSPG>

Generative Adversarial Networks: Computer Science Meets Biology

Generative adversarial networks (GANs) are deep neural net architectures comprised of two nets, pitting one against the other (thus the “adversarial”). “GANs’ potential is huge, because they can learn to mimic any distribution of data.”

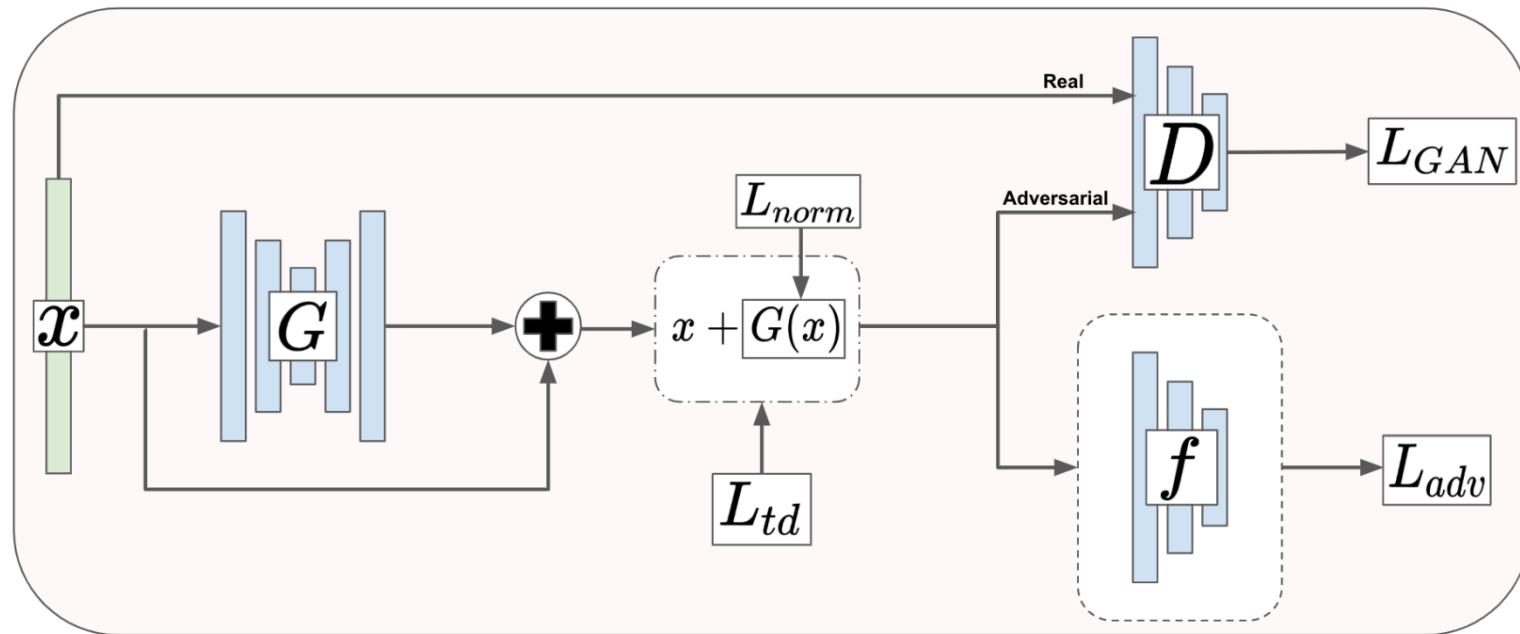


<https://codewords.recurse.com/issues/five/why-do-neural-networks-think-a-panda-is-a-vulture>

<https://skymind.ai/wiki/generative-adversarial-network-gan>

The Transcriptome State Perturbation Generator (TSPG)

Goal: Learn to produce x_{adv} classified as $f(x_{adv}) = t$ where t is a target class (e.g. tissue A) different from ground truth label y (e.g. tissue B).



G=generator network

D=discriminator

X=features (genes)

f=pretrained target network

$G(x)$ =feature perturbation

$X_{adv} = X + G(x) = \text{adversarial example}$

$L = L_{GAN} + L_{adv} + L_{norm} + L_{td} = \text{Loss for } G$

Accuracy of Target Network on Perturbed GTEx Datasets

Gene Set	Genes	Target Class	f Accuracy (%)
Hallmark Hedgehog Signaling	36	nerve-tibial	100
Hallmark Peroxisome	107	brain-spinal cord	100
Hallmark Apoptosis	161	lung	99.9
Hallmark E2F Targets	200	artery-coronary	99.8
Hallmark All	4386	thyroid	100
Hallmark All	4386	heart-left ventricle	100

Tricking the Hippocampus to be a Tibial Nerve

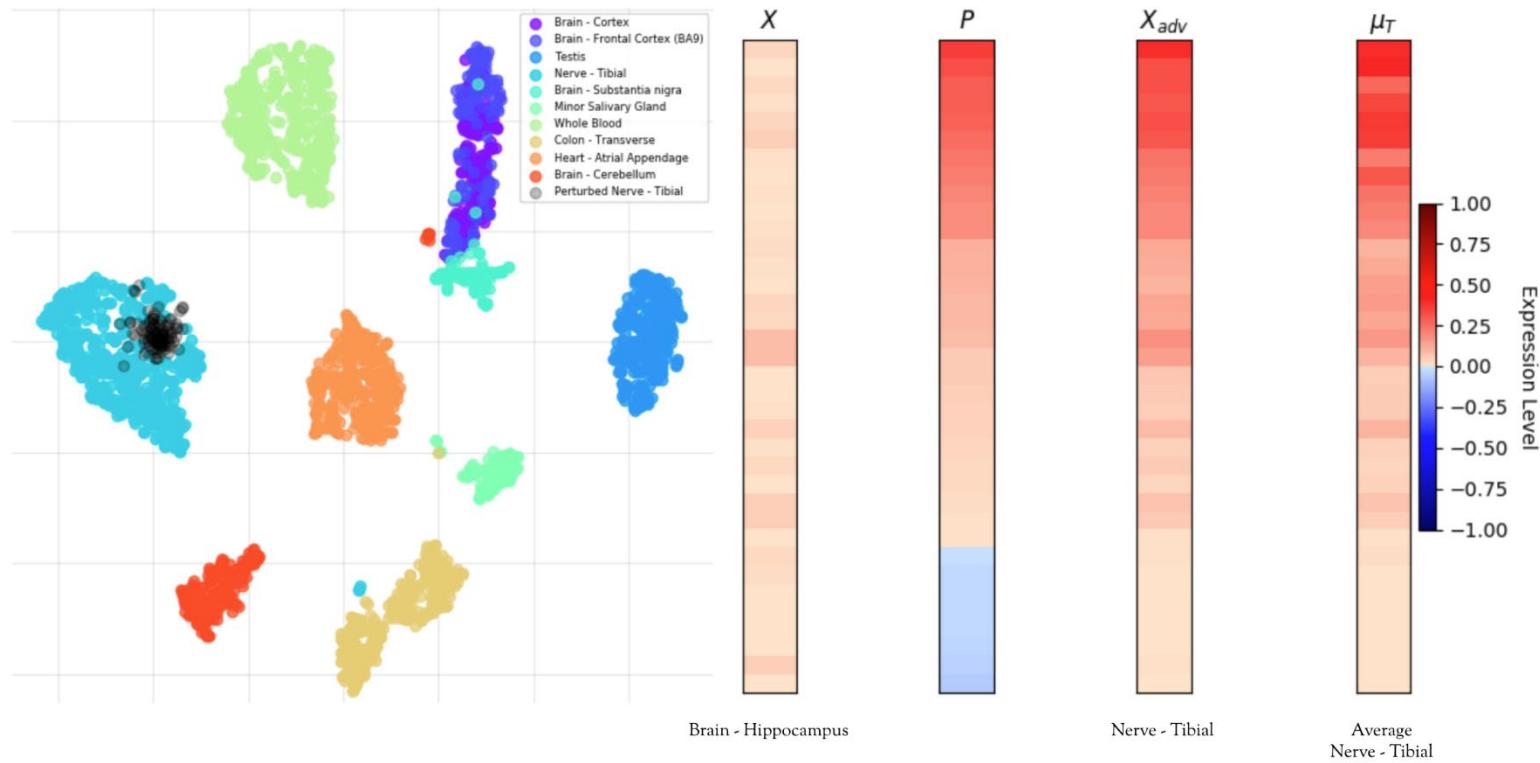


Figure 1: Adversarial generation for Nerve-Tibial target using the Hallmark Hedgehog Signaling gene set. t-SNE plot of original and perturbed samples using the Hallmark Hedgehog Signaling gene set (left). Heatmap of cellular transformation from Brain-Hippocampus to Nerve-Tibial (right). Perturbation (P) ranges from $[-1, 1]$, which is added to original sample (X), then adversarial example (x_{adv}) is clipped to $[0, 1]$. The mean expression vector (μ_T) of the target class (Nerve-Tibial) is shown.

Patient-Specific Tumor Analysis

 **CellPress**
OPEN ACCESS

Patterns
Article

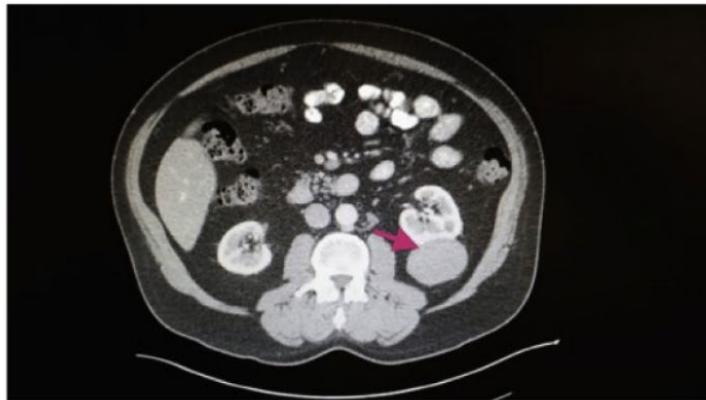


Figure 4. Computed Tomography Image of Patient BP's Renal Cell Carcinoma
Imaged in February of 2014, this thoracic computed tomogram with intravenous contrast shows a 5.7 × 4.8-cm mass in the left kidney.



Table 3. Sample Counts Included in the Comprehensive Kidney Cancer GEM

Tissue Type	Dataset of Origin	Count
KIRC tumor	TCGA	475
KIRC normal	TCGA	72
KIRP tumor	TCGA	236
KIRP normal	TCGA	29
KICH tumor	TCGA	60
KICH normal	TCGA	25
Kidney normal	GTEX	32
BP tumor	-	1
BP normal	-	1

Patterns

Volume 1, Issue 6, 11 September 2020, 100087



Article

Cellular State Transformations Using Deep Learning for Precision Medicine Applications

Colin Targonski ^{1, 6}, M. Reed Bender ^{2, 6}, Benjamin T. Shealy ¹, Benafsh Husain ¹, Bill Paseman ³, Melissa C. Smith ¹, F. Alex Feltus ^{2, 4, 5, 7, 8, 9}

RareKidneyCancer.org

TSPG Discovers Kidney Cancer Biomarkers

Patterns
Article

CellPress
OPEN ACCESS

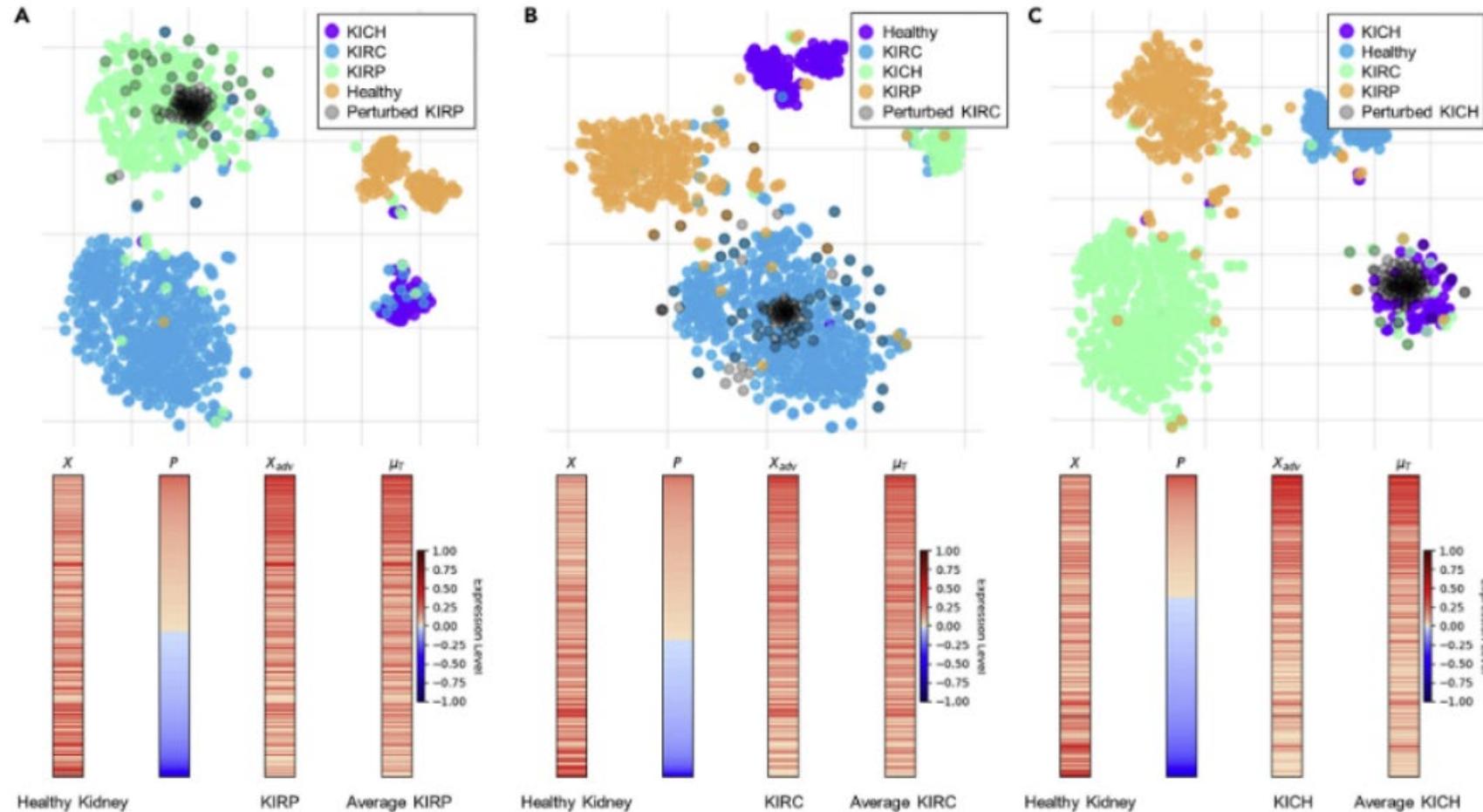


Figure 3. Adversarial Generation for Three Subtypes of Kidney Cancer Using All Hallmark Genes as the Input Gene Set

t-SNE plot and corresponding heatmap of cellular transformation from healthy to KIRP (A), healthy to KIRC (B), and healthy to KICH (C). Perturbations (P) range from $[-1, 1]$, which is added to original sample (X), then adversarial sample (X_{adv}) is clipped to $[0, 1]$. The mean expression vector (μ_T) of the target class is shown.

Geographically Distributed Interdisciplinary Science is Super Fun!

Feltus Lab

Yueyao Gao (<PhD, G&B)
Cole McKnight (Cloud Architect)
Rini Pauly (<PhD, G&B)
Kaitlyn Gmitro (<BSc, G&B)
Meghan Mobley (<BSc, G&B)
Brad Selee (<BSc, ECE)
Devin Keane (<MSc, G&B)

Recent alumni

Nicole Nelligan (BSc, G&B)
Ben Shealy (PhD, ECE)
Allison Hickman (PhD, G&B)
Yuqing Hang (PhD, G&B)
Reed Bender (MSc, BioE)
Ethan Bensman (BSc, CS)
Cameron Ogle (BSc, CS)
Benafsh Husain (PhD, BDSI)
Will Poehlman (PhD, G&B)
Jordan Little (BSc, G&B)
Rachel Eimen (Bsc, ECE)
Courtney Shearer (BSc, CS)
Melissa Judge (BSc, BioE)
Colin Targonski (MSc, , ECE)
Leland Dunwoodie (<MD, G&B)
Olivia Feltus (<BSc, Intern)
Nick Watts (Software Engineer)
Zach Gerstner (MS, BDSI)
Cole Younginer (<MSc, ECE)



@ Clemson

Melissa Smith (ECE)
KC Wang (ECE)
Walt Ligon (ECE)
Nick Mills (ECE)
Jon Calhoun (ECE)
Brian Dean (CS)
Marc Birtwistle (ChemE)
Brian Booth (BioE)
Julia Frugoli (G&B)
Suchitra Chavan (G&B)
Elise Schnabel (G&B)
Susan Duckett (AVS)
Corey Ferrier (CCIT)
Jim Pepin (CCIT)
Xizhou Feng (CCIT)

@ Earth

Stephen Ficklin (WSU)
Josh Burns (WSU)
Tyler Biggs (WSU)
Dorrie Main (WSU)
Sook Jung (WSU)
Joe Breen (Utah)
Jill Wegrzyn (UCONN)
Meg Staton (UTK)
Jim Bottum (Internet2)
Dana Brunson (Internet2)
John Hicks (Internet2)
Marvin Weinstein(QI) Ken
Matusow (Synergy)
Karan Sapra (nVidia)
David Clarke (Toolwire)
Jacob Loftis (Cisco)
Mike Shepherd (Cisco)
Mike Kowal (Cisco)
Jordan Auge (Cisco)

@ Earth (cont.)

Mats Rynge (USC-OSG)
Bala Desinghu (U Chicago-OSG)
Andrew Paterson (UGA)
Claris Castillo (RENCI)
Steve Cox (RENCI)
Ray Idaszak (RENCI)
Paul Ruth (RENCI)
Michael Stealy (RENCI)
Isma Gilani (RENCI)
Fan Jiang (RENCI)
Mert Cevik (RENCI)
Ananya Mukherghee (RENCI)
Emily Casanova (USC-GHS)
Manual Casanova (USC-GHS)
Alex Bowers (Columbia U.)
Josh Vandenbrink (Louisiana Tech)
Ann Loraine (UNCC)
Colleen Doherty (NCSU)
John Graham (UCSD)
Wallace Chase (REANNZ)
Christos Papadopoulos (Memphis)
Susmit Shannigrahi (Tennessee Tech)
Chengyu Fan (CSU)

Many many more!



Thank You Funding Agencies!!!!

- “*CC* Integration-Large: Prototyping a Secure Distributed Storage Infrastructure for Accelerating Big Science.*” [2126148] (S. Shannigrahi PI)
- “*CC*Data: National Cyberinfrastructure for Scientific Data Analysis at Scale (SciDAS)* Source: NSF-CC* [1659300] (A. Feltus PI)
- “*RCN: Advancing Research and Education Through a National Network of Campus Research Computing Infrastructures - The CaRC Consortium*” Source: NSF [1620695] (J Bottum PI > A. Feltus PI)
- “*MRI: Acquisition of a Cyberinstrument for Interdisciplinary Computational Science and Engineering.*” Source: NSF-MRI [1725573] (A. Apon PI)



Historical

- “*MCA-PGR: Spatial and Temporal Resolution of mRNA Profiles During Early Nodule Development.*” Source: NSF-PGRP [1444461] (J. Frugoli PI)
- “*Tripal Gateway: Platform for Next-Generation Data Analysis and Sharing.*” Source: NSF-DIBBS [1443040] (S. Ficklin, PI)
- “*BIGDATA: F: DKM: Collaborative Research: PXFS: ParalleX Based Transformative I/O System for Big Data*” Source: NSF-BIGDATA [1447771] (W. Ligon PI)
- “*Genomic and Breeding Foundations for Bioenergy Sorghum Hybrids.*” Source: Plant Feedstock Genomics for Bioenergy [DE-FOA-000041] (S Kresovich, PI).
- “*Big Data Visualization REU*”. Source: National Science Foundation [1359223](V Byrd, PI)
- “*MRI: Acquisition of a High Performance Computing Instrument for Collaborative Data-Enabled Science.*” Source: National Science Foundation [1228312] (A Apon, PI)
- “*CC-NIE Integration: Clemson-NextNet*” Source: National Science Foundation [1245936] (KC Wang, PI)
- “*Building non-model species genome curation communities.*” Source: National Evolutionary Synthesis Center (NESCent) (A Papanicolaou, PI)
- “*Big Data Analysis Tools for Agricultural Genomics.*” Source: Clemson University Experiment Station (USDA Hatch Project) [SC-1700492] (A. Feltus PI).

