

BIOE 310:
Computational Tools for
Biological Data

What this class is all about?

Instructor

- Name: **Sergei Maslov**
- **Professor of Bioengineering, Physics, Carl R. Woese Institute for Genomic Biology, and National Center for Supercomputing Applications**
- Office: 3103 Carl Woese Institute for Genomic Biology and sometimes 3146C Everitt Laboratory (both by appointment)
- E-mail: maslov@illinois.edu
- Phone: 217-265-5705



Teaching Assistant:

Seokjin Yeo

sy44@illinois.edu



Questions and Suggestions:

maslov@illinois.edu

sy44@illinois.edu

Start subject with [BIOE310]

Homework and Exams

- **Homework assignments.**

Due at the beginning of the class on the designated day

- **Midterm exam.** March either before or after the spring break

- **Final exam.** Date will be decided by the College of Engineering

- **Grading:**

Homework	30%
Midterm	30%
Final	40%

Course Website

<https://courses.engr.illinois.edu/bioe310>

Grades will be on

<https://my.bioen.illinois.edu/gradebook>

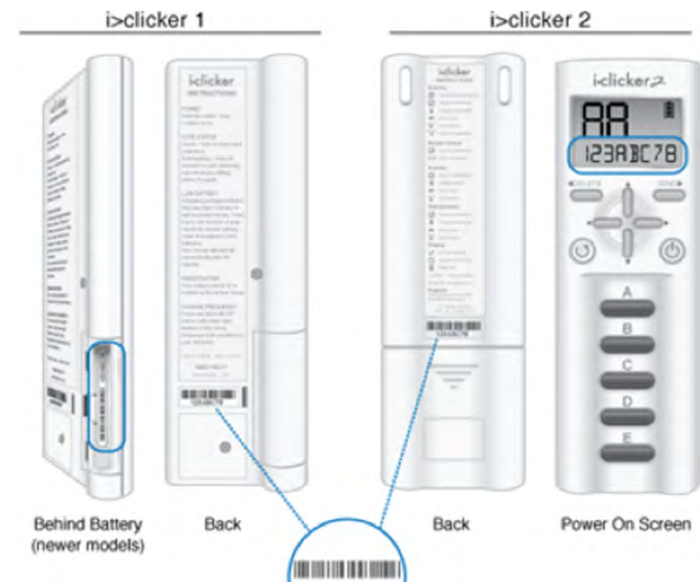
BIOE 310 - Computational Tools for Biological Data

[Return to syllabus](#)

#	Date	Topics	Slides	Matlab	Homework	Exams
1	Jan 26					
2	Jan 28					

Bring your iClickers to my lectures

- Who knows what is an iClicker?
- Show of hands: who has an iClicker?
- I would like you all to have an iClicker and bring it to every class. On **amazon.com** a new **iClicker** (1st generation is OK) costs around \$40. It is also sold at UIUC Bookstore. The used ones are cheaper.
- An alternative solution is using a mobile app:
<https://www.iclicker.com/students/apps-and-remotes/apps>
- Your answers **WILL NOT** be used for grading.
I need them to see if I lost some of you and what could I rephrase to better explain the material



We will use Matlab in class

- Bring **your laptops to class**
- **Poll: who has Matlab?**
- Need to have **Matlab installed** and know the basic user interface (inline commands, plotting)
- We will use **Statistics and Machine Learning Toolbox and Bioinformatics Toolboxes**
- You can use CITRIX for UIUC students and connect to EWS Windows Lab Software
- **.m files and .mat** with Matlab commands and data **will be on the website** after the lecture

Who has Matlab?

- A. Have on my own laptop
- B. Plan to use CITRIX
- C. I don't have Matlab
- D. I don't know yet
- E. I will never use Matlab

Get your i-clickers

We will use Matlab in class

- Bring **your laptops to class**
- Need to have **Matlab installed** and know the basic user interface (inline commands, plotting)
- We will use **Statistics and Machine Learning Toolbox and Bioinformatics Toolboxes**
- Good news! Now all faculty and graduate students get Matlab **for free**. See [offering on the WebStore](#) site and follow the [detailed instructions](#).
- **.m files and .mat** with Matlab commands and data **will be on the website** after the lecture

Possible alternative to purchasing Matlab and toolboxes is to use campus resources.

Both Engineering Workstations (EWS) and ACES computers have Matlab.
I don't think all of them offer the statistics and bioinformatics toolboxes
(EWS should, ACES computers may not..).

See the following to access:

Citrix for EWS, Matlab, and ACES computers -- links for all

<https://it.engineering.illinois.edu/ews/lab-information/remote-connections/connecting-citrix>

<https://it.engineering.illinois.edu/services/instructional-services/remote-connections-citrix>

Accessing Engineering Workstations (EWS)

<https://it.engineering.illinois.edu/ews>

Accessing ACES Academic Computing Workstations

<http://acf.aces.illinois.edu/remote/>

<http://acf.aces.illinois.edu/remote/pc.html>

To access off campus use:

CISCO Virtual Private Network -- For off-campus access to campus computer and network resources
(software programs, files saved on the network, etc.)

<https://techservices.illinois.edu/services/virtual-private-networking-vpn/download-and-set-up-the-vpn-client>

CISCO VPN CLIENT

<https://webstore.illinois.edu/shop/product.aspx?zpid=2600>

CISCO AnyConnect VPN

<https://webstore.illinois.edu/shop/product.aspx?zpid=1222>

What will you learn in this course?

- Basics of probability and statistics
 - Basic concepts of probability, Bayes theorem
 - Discrete and continuous probability distributions
 - Multivariate statistics
 - Sampling distributions
 - Parameter estimation
 - Hypothesis testing
 - Regression
- How it is applied to biological data
 - Basics of genomics
 - Systems biology (gene expression, networks)

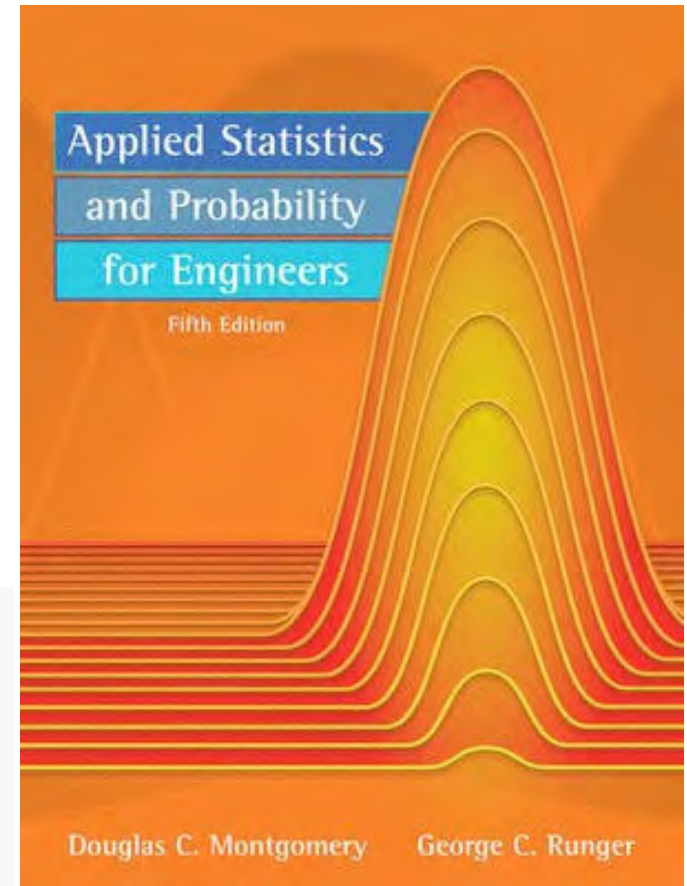
The main Probability/Statistics Textbook

Applied Statistics and Probability for Engineers, 5th Edition

D. C. Montgomery and G. C. Runger
John Wiley & Sons, Inc. (2011)

You can also use other editions from
4th (2007) to 6th (2014)

5th edition is available for free
at our library



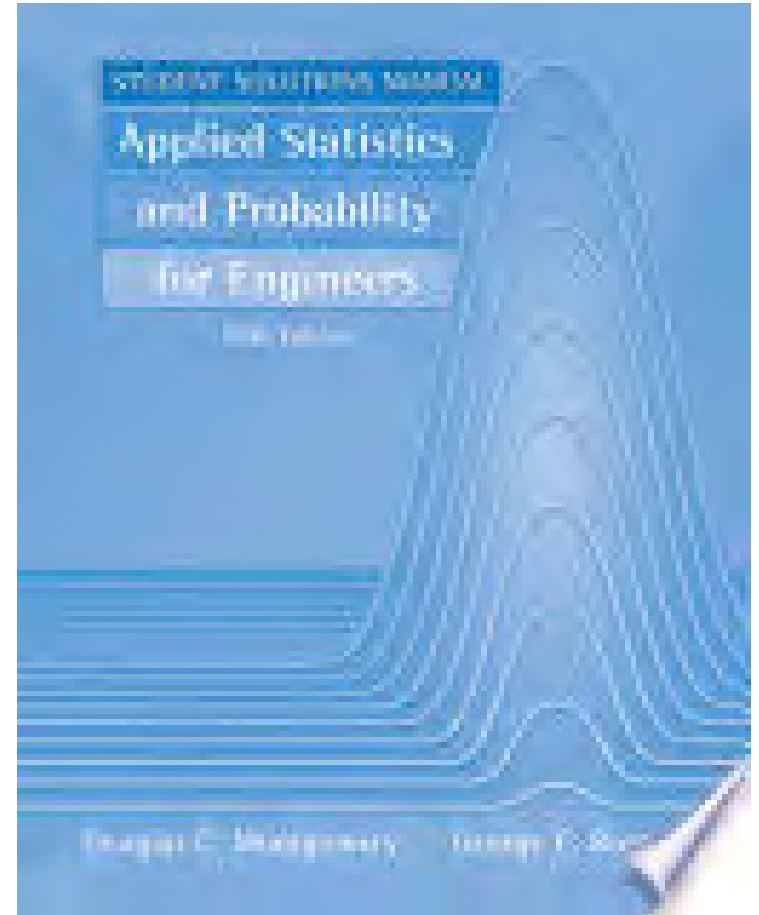
Problems for our main Probability/Statistics Textbook

Student Solutions Manual Applied Statistics and Probability for Engineers, 5th Edition

D. C. Montgomery and G. C. Runger
John Wiley & Sons, Inc. (2010)

You can also use other editions from
4th (2007) to 6th (2014)

5th edition is available
for free at our library



Probability/Statistics for Bioengineering with Matlab exercises

Statistics for Bioengineering Sciences

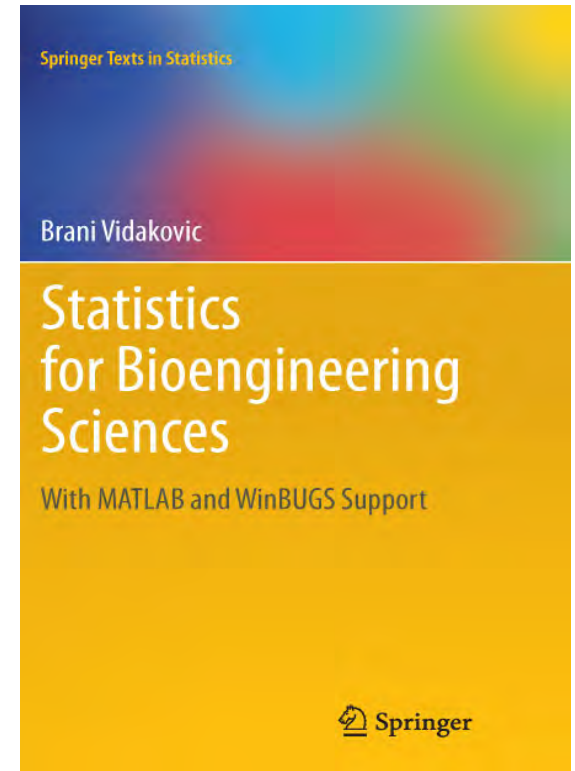
with MATLAB and WinBUGS Support

Brani Vidakovic

Department of Biomedical Engineering, Georgia Tech

(2011) Springer, New York

*It is constantly updated with the newest version at the link
below.*



Free as a PDF eBook at

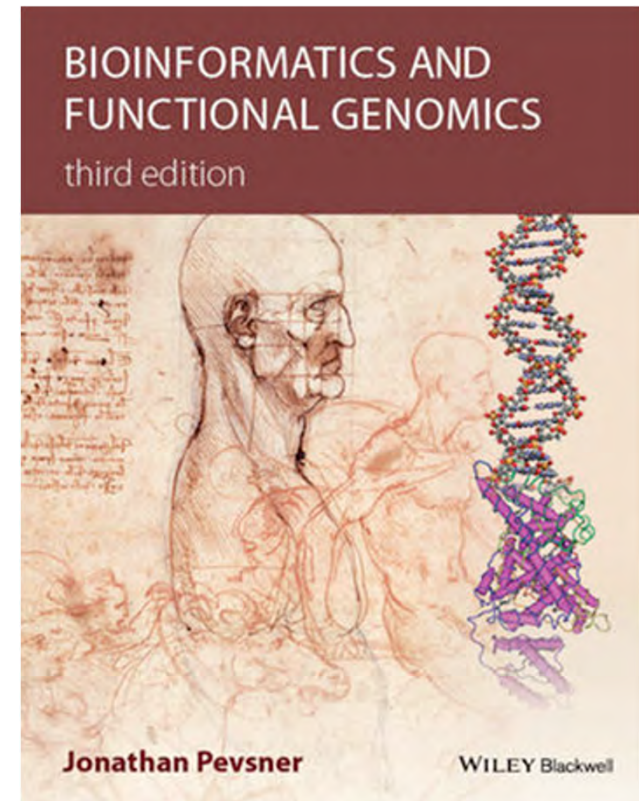
<http://statbook.gatech.edu/statb4.pdf>

Matlab exercises and datasets are at

<http://springer.bme.gatech.edu>

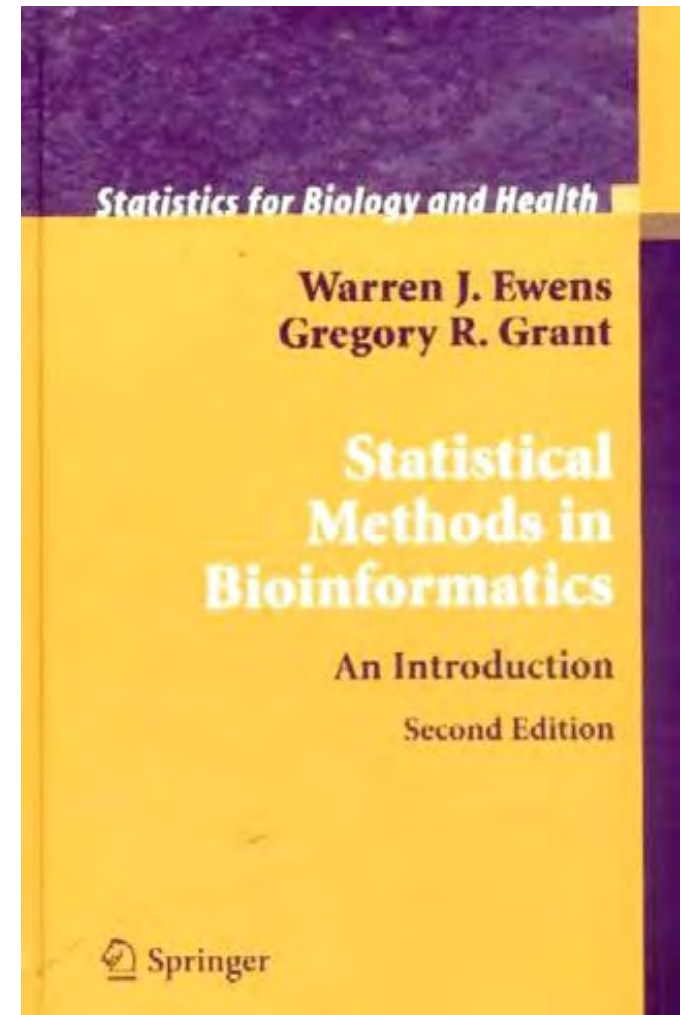
Genomics/Systems Biology Textbook

- J Pevsner
Bioinformatics and functional genomics
Wiley-Blackwell,
2nd edition [2009] exists in electronic form
3rd edition [2015] *has up-to-date
information on NGS: RECOMMENDED
(about \$60 on amazon)*
- 2nd edition is available for free
in electronic form in our library



Another Bioinformatics/Statistics Textbook

- *Ewens, WJ and Grant, GR Statistical Methods in Bioinformatics: An Introduction, 2nd ed, Springer, 2005.*
- *2nd edition as PDF eBook*



Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE SO MANY CROWS IN ROCHESTER, MN

WHY IS PSYCHIC WEAK TO BUG

WHY DO CHILDREN GET CANCER

WHY IS POSEIDON ANGRY WITH ODYSSEUS

WHY IS THERE ICE IN SPACE

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED

WHY IS SPACE BLACK

WHY IS OUTER SPACE SO COLD

WHY ARE THERE PYRAMIDS ON THE MOON

WHY IS NASA SHUTTING DOWN

WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE
GHOSTS



WHY IS THERE AN OWL IN MY BACKYARD

WHY IS THERE AN OWL OUTSIDE MY WINDOW

WHY IS THERE AN OWL ON THE DOLLAR BILL

WHY DO OWLS ATTACK PEOPLE

WHY ARE AK 47s SO EXPENSIVE

WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE

WHY ARE THERE GODS

WHY ARE THERE TWO SPOCKS

WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS

WHY AREN'T MY QUAIL EGGS HATCHING

WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE

WHY AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS

WHY DO I SAY UH
WHY IS SEA SALT BETTER
WHY ARE THERE TREES IN THE MIDDLE OF FIELDS

WHY IS THERE NOT A POKEMON MMO
WHY IS THERE LAUGHING IN TV SHOWS
WHY ARE THERE DOORS ON THE FREEWAY

WHY ARE THERE SO MANY SVCHOST.EXE RUNNING
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA
WHY ARE THERE SCARY SOUNDS IN MINECRAFT

WHY IS THERE KICKING IN MY STOMACH
WHY ARE THERE TWO SLASHES AFTER HTTP
WHY ARE THERE CELEBRITIES

WHY DO SNAKES EXIST
WHY DO OYSTERS HAVE PEARLS
WHY ARE DUCKS CALLED DUCKS

WHY DO THEY CALL IT THE CLAP
WHY ARE KYLE AND CARTMAN FRIENDS
WHY IS THERE AN ARROW ON AANG'S HEAD

WHY ARE TEXT MESSAGES BLUE
WHY ARE THERE MUSTACHES ON CLOTHES
WHY ARE THERE MUSTACHES ON CARS

WHY ARE THERE MUSTACHES EVERYWHERE
WHY ARE THERE SO MANY BIRDS IN OHIO
WHY IS THERE SO MUCH RAIN IN OHIO

WHY IS OHIO WEATHER SO WEIRD
WHY ARE THERE MALE AND FEMALE BIKES
WHY ARE THERE BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP
WHY AREN'T THERE VARICOSE ARTERIES
WHY ARE OLD KINGDOMS DIFFERENT

WHY ARE THERE SQUIRRELS
WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR

WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE

WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP

WHY DO DREAMS SEEM SO REAL

WHY ARE THERE SO MANY SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

This course is about biological data
and probability theory, and statistics
concepts needed for its analysis

What biological data will be discussed?

Will be covered in lectures or Matlab exercises:

- Genomic data: strings of letters ACGT
- Gene Expression data: messenger RNA copy numbers transcribed from genes
- Proteomic data: protein abundances
- Network data: pairs of interacting genes or proteins and protein-protein interaction strengths

Will not be covered:

- Imaging data such as e.g. fMRI brain scans, Brain connectome data, Ecosystem dynamics data

Why do you need
probability and statistics
to analyze
modern biological data?

Definition of **probability theory** by Encyclopedia Britannica

a branch of mathematics concerned
with the analysis of **random
phenomena**

Definition of ***statistics*** by Merriam-Webster

1 : a branch of mathematics dealing with the
collection, analysis, interpretation, and
presentation of **masses of numerical data**

...

Why do you need
probability and statistics
to analyze
modern biological data?

Reason 1:
Biology now has Lots of Data

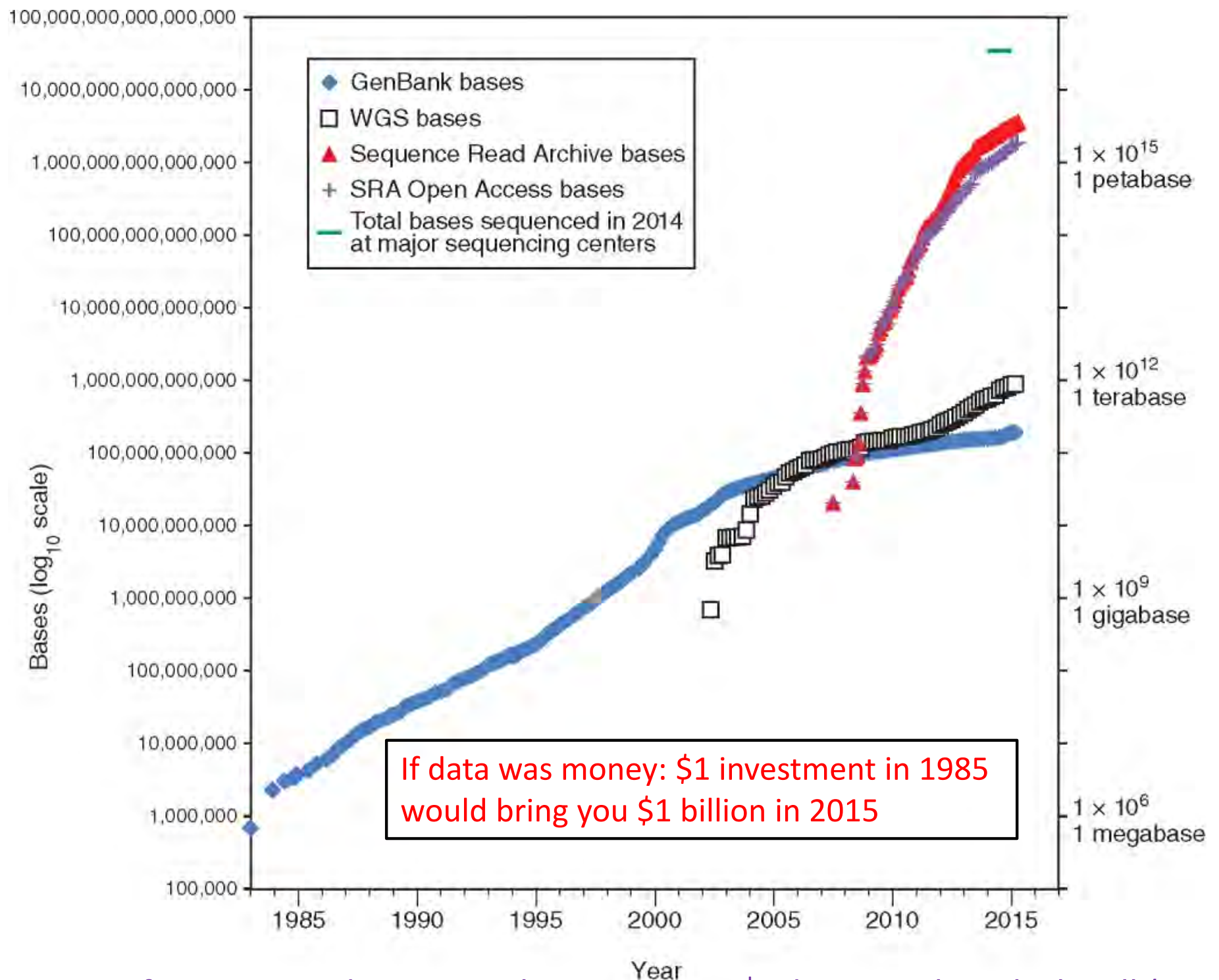
Base pairs	Unit	Abbreviation	Example
1	1 base pair	1 bp	A, C, G, T = 2 bits = 0.25 bytes
1000	1 kilobase pair	1 kb	
1,000,000	1 megabase pair	1 Mb	
10^9	1 gigabase pair	1 Gb	
10^{12}	1 terabase pair	1 Tb	
10^{15}	1 petabase pair	1 Pb	

3

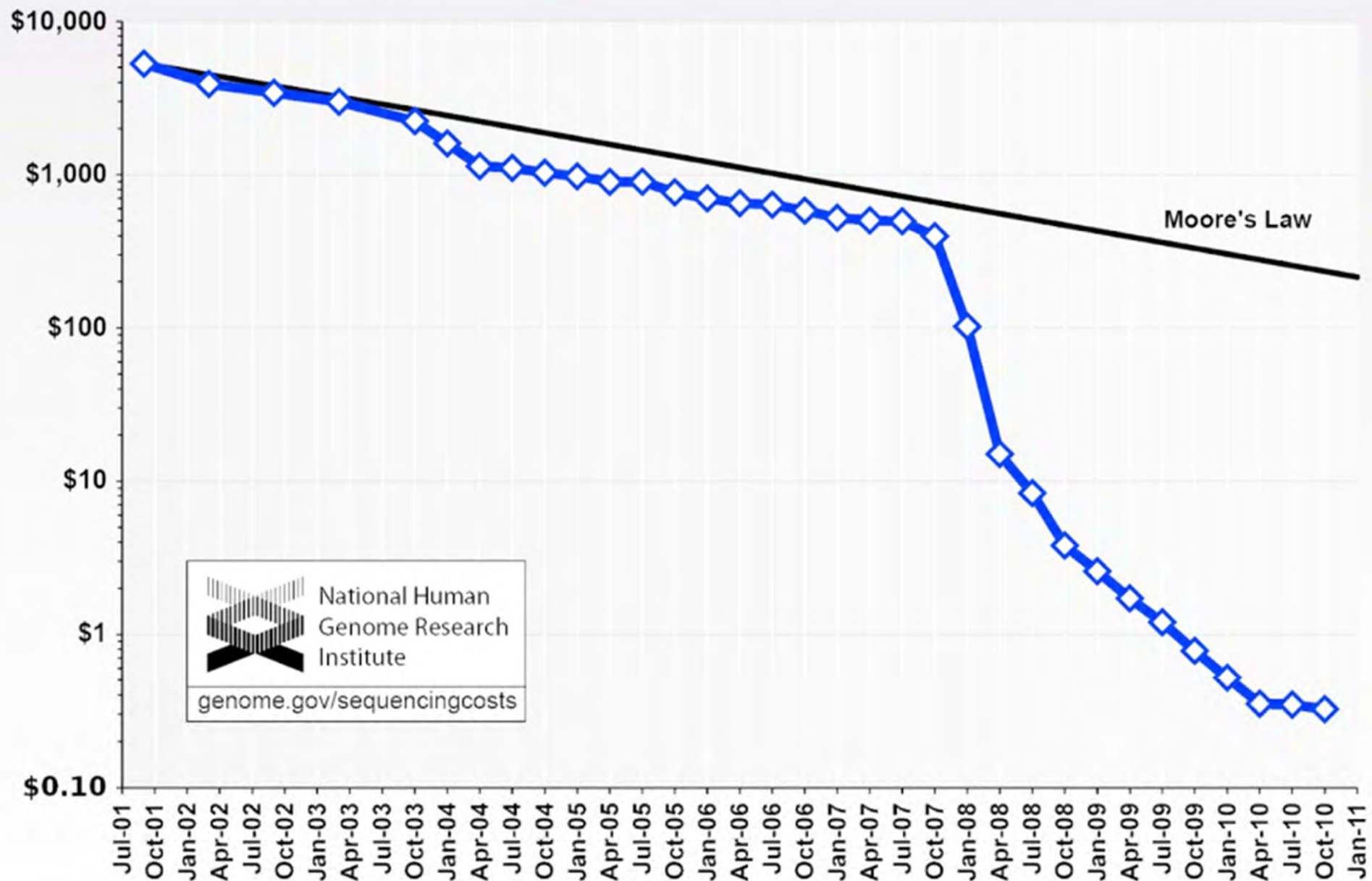
Size	Abbreviation	No. bytes	Examples
Bytes	–	1	1 byte is typically 8 bits, used to encode a single character of text
Kilobytes	1 kb	10^3	Size of a text file with up to 1000 characters
Megabytes	1 MB	10^6	Size of a text file with 1 million characters
Gigabytes	1 GB	10^9	600 GB: size of GenBank (uncompressed flat files) ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt (WebLink 2.84)
Terabytes	1 TB	10^{12}	385 TB: <u>United States Library of Congress web archive</u> (http://www.loc.gov/webarchiving/faq.html) (WebLink 2.85) 464 TB: Data generated by the <u>1000 Genomes Project</u> (http://www.1000genomes.org/faq/how-much-disk-space-used-1000-genomes-project) (WebLink 2.86)
Petabytes	1 PB	10^{15}	1 PB: size of dataset available from <u>The Cancer Genome Atlas (TCGA)</u> 5 PB: size of <u>SRA data available for download from NCBI</u> 15 PB: amount of data produced <u>each year at the physics facility CERN (near Geneva)</u> (http://home.web.cern.ch/about/computing) (WebLink 2.87)
Exabytes	1 EB	10^{18}	<u>2.5 exabytes of data are produced worldwide (Lampitt, 2014)</u>

Bacterial genome = "War & Peace"

Human Genome

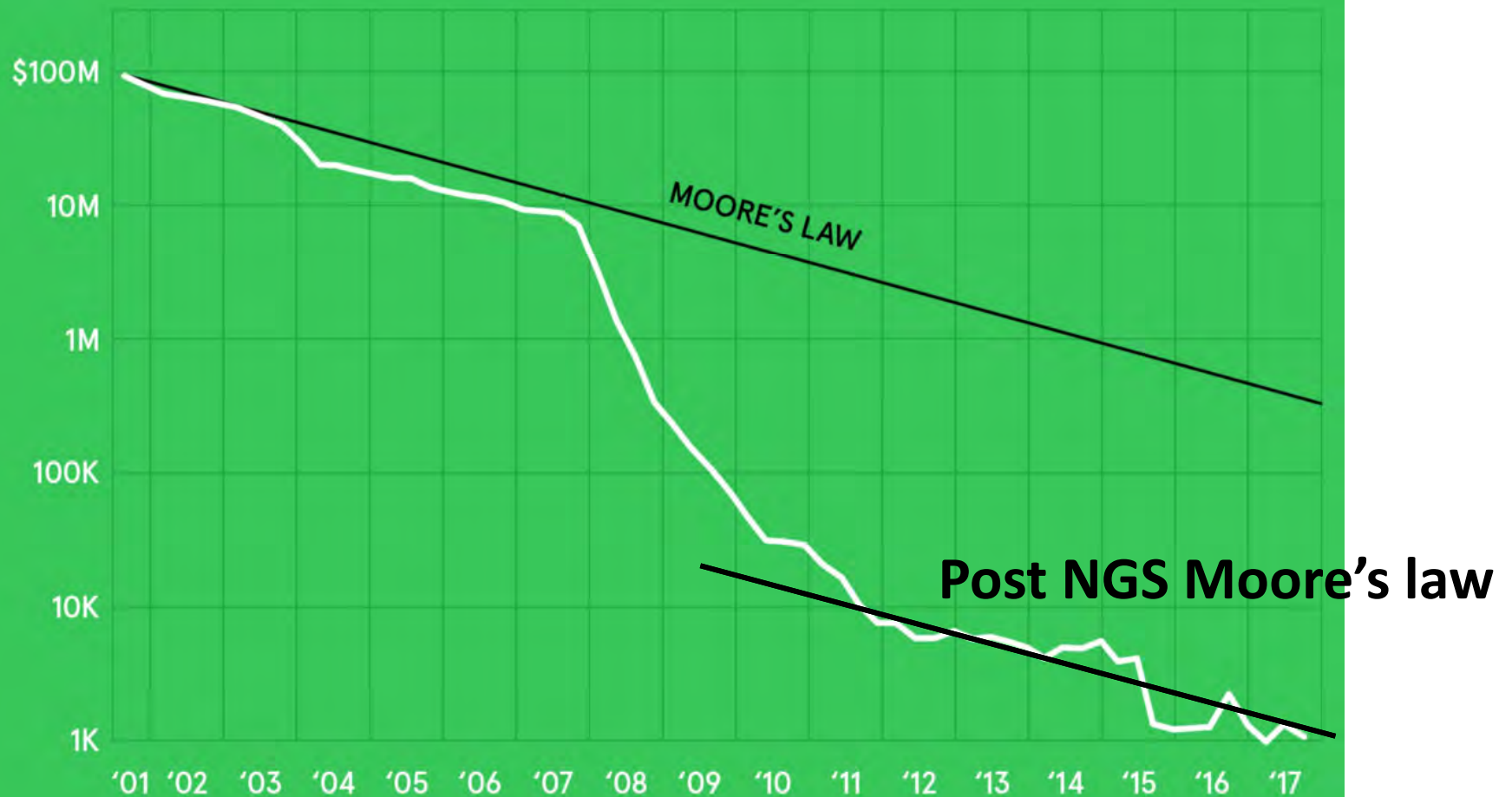


Cost per Megabase of DNA Sequence



Cost per Genome Sequenced

The cost of sequencing a human genome compared with the reductions that would be expected at the rate Moore's law predicts for computer chips. Over the past decade, next-generation sequencing and cloud computing drove the figure down. The average bumped higher in recent years because of brief slowdowns in production.



Source: NIH

NEO LIFE

Who will have **bigger data** by 2025?

<u>Data Phase</u>	<u>Astronomy</u>	<u>Twitter</u>
Acquisition	25 zetta-bytes/year	0.5–15 billion tweets/year
Storage	1 EB/year	1–17 PB/year
<div>Peta=10^{15} Exa=10^{18} Zetta=10^{21}</div>		
<u>YouTube</u>	<u>Genomics</u>	
500–900 million hours/year	1 zetta-bases/year	
1–2 EB/year	2–40 EB/year	

Z. Stephens, S. Lee, F. Faghri, R. Campbell, C. Zhai, M. Efron,
R. Iyer, M. Schatz, S. Sinha, and G. Robinson (2015) PLoS Biol 13: e1002195.

Growth of DNA Sequencing

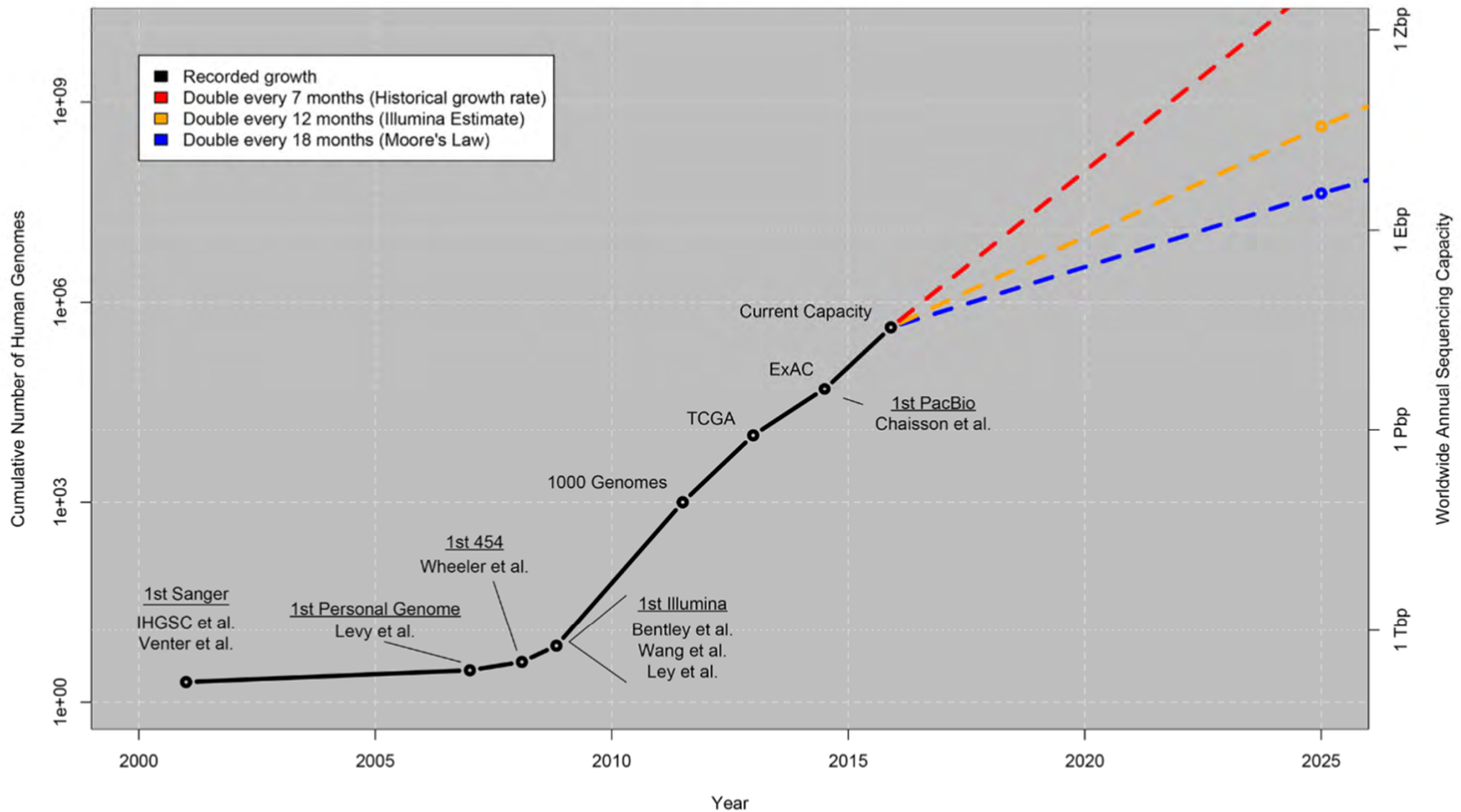
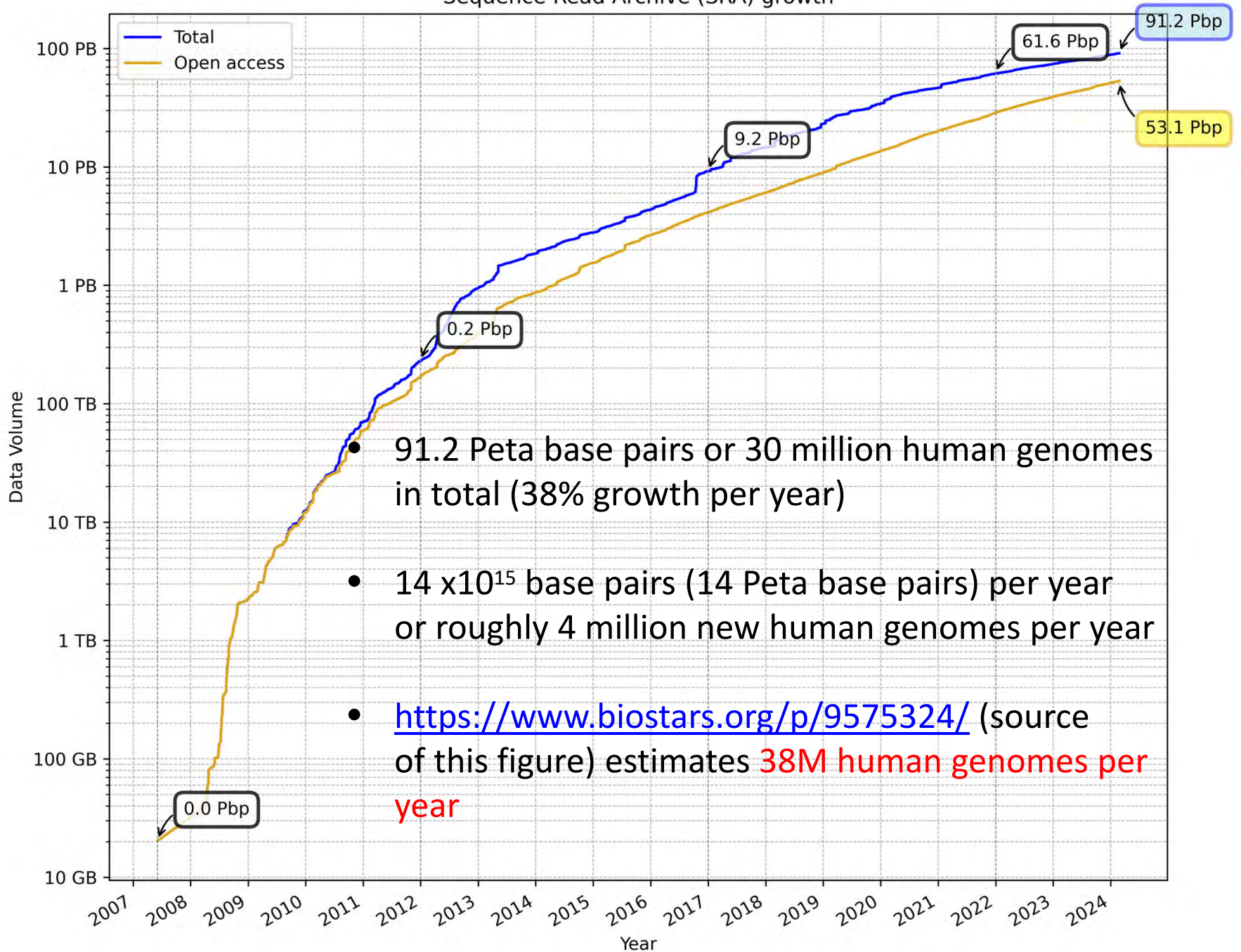


Fig 1. Growth of DNA sequencing. The plot shows the growth of DNA sequencing both in the total number of human genomes sequenced (left axis) as well as the worldwide annual sequencing capacity (right axis: Tera-basepairs (Tbp), Peta-basepairs (Pbp), Exa-basepairs (Ebp), Zetta-basepairs (Zbps)). The values through 2015 are based on the historical publication record, with selected milestones in sequencing (first Sanger through first PacBio human genome published) as well as three exemplar projects using large-scale sequencing: the 1000 Genomes Project, aggregating hundreds of human genomes by 2012 [3]; The Cancer Genome Atlas (TCGA), aggregating over several thousand tumor/normal genome pairs [4]; and the Exome Aggregation Consortium (ExAC), aggregating over 60,000 human exomes [5]. Many of the genomes sequenced to date have been whole exome rather than whole genome, but we expect the ratio to be increasingly favored towards whole genome in the future. The values beyond 2015 represent our projection under three possible growth curves as described in the main text.

Sequence Read Archive (SRA) growth



<https://www.ncbi.nlm.nih.gov/sra/docs/sragrowth/>

What makes genomic data so big?

- There are **~9 millions species** each with its own genome
- **Each of us humans** (7.5 billions and counting) has **unique DNA**: we want to compare them all to each other
- Each cell has **just 1 genome (DNA)** but **multitude of transcriptomes (RNA levels)** and **proteomes (protein levels)**
- **Cancer cells acquire mutations** in their genomes: need to track **multiple lineages in a tumor vs time** to understand cancer
- **DNA** was proposed as a **long-term storage medium** of information

Farfetched? Storage standards evolve fast but DNA standard remained unchanged for 4 billion years

Note: Nature article started the comparison with a hard drive and flash memory skipping the floppy disk



How DNA could store all the world's data

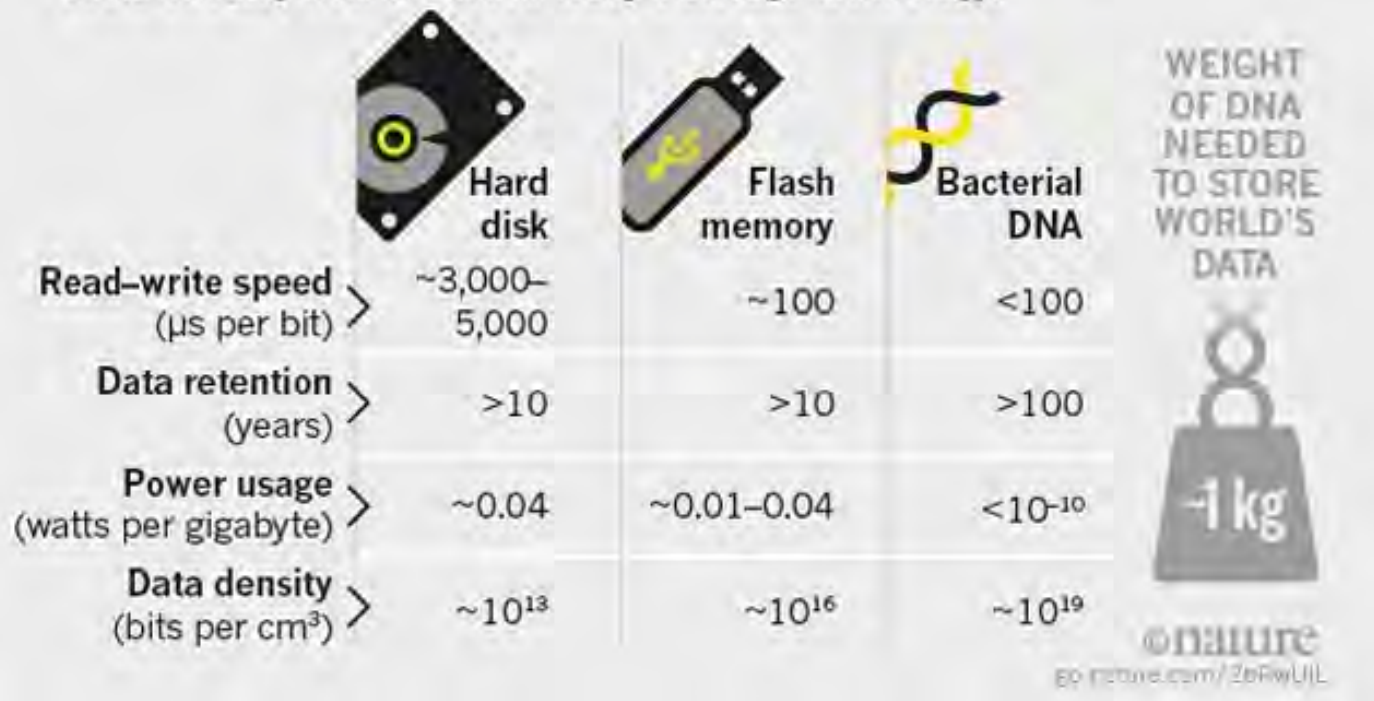
Modern archiving technology may hold an answer to that problem

Andy Extnance

31 August 2016

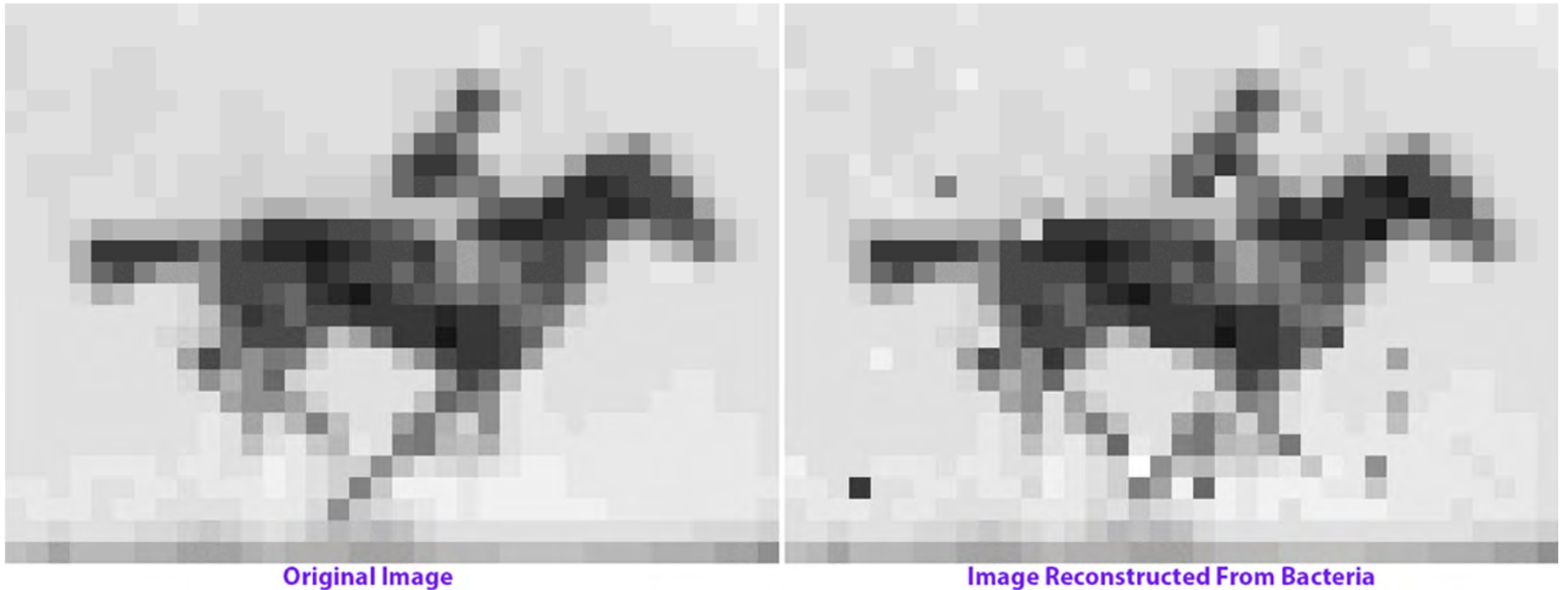
STORAGE LIMITS

Estimates based on bacterial genetics suggest that digital DNA could one day rival or exceed today's storage technology.



- Prof Olgica Milenkovic from Electrical and Computer Engineering UIUC is a local expert on this topic
- Profs. George Church and Sri Kosuri (Harvard Medical School) explains a potential use of DNA as storage medium in 2012
- <https://www.youtube.com/watch?v=IJAdqAVjQqY>

Fast-forward from 2012 to 2017



Shipman SL, Nivala J, Macklis JD, Church GM.
CRISPR–Cas encoding of a digital movie into the genomes
of a population of living bacteria. *Nature*. 2017;547: 345–349. doi:10.1038/nature23017

Why do you need
probability and statistics
to analyze
modern biological data?

Reason 2:
Life is random and messy

Show video “Cell organelles”

- Made at the Walter and Eliza Hall Institute of Medical Research at Victoria, Australia
- Animated by award-winning artist Dr. Drew Berry
- Go to <https://www.wehi.edu.au/wehi-tv> for other videos

Life is messy, random, and noisy

Yet it is beautifully complex
and has many parts
(see statistics)

Why life is so random?

- Biomolecules are very small
(nano- to micro-meters) → Brownian noise
- # molecules/cell is often small →
Large cell-to-cell variations
- Genomic data comes from biological evolution
– the Mother of all random processes
- Genomic data involves (random) samples
– We have genomes of some (not all) organisms
– We have tissue samples of some (not all) cancer patients

Why life is so complex?

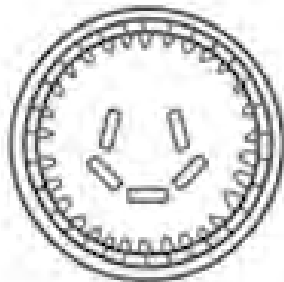
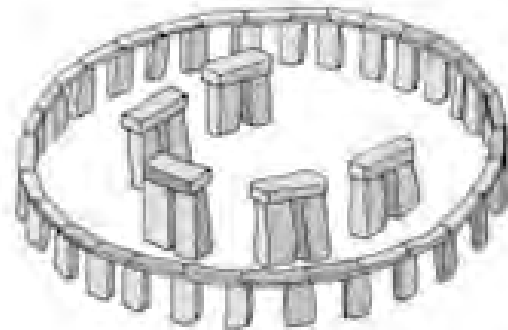
Primer on complex system

Complex systems have many interacting parts

- All **parts** are **different** from each other
 - 10s thousands (10^4) types of **proteins** in an organism
 - 100 thousands (10^5) **organizations (AS)** in the Internet
 - 1 billion (10^9) people on **Facebook**
 - 10 billion (10^{10}) **web pages** in the WWW
 - 100 billion (10^{11}) **neurons** in a human brain
 - **NOT 10^{23} electrons or quarks studied by physics: they are all the same and boring!**
- Yet they **share** the same **basic design**
 - All proteins are strings of the **same 20 amino acids**
 - All WWW pages use **HTML**, JavaScript, etc.
 - All neurons generate and receive **electric spikes**

Example: a complex system with many parts

HËNJ



80x



30x



30x



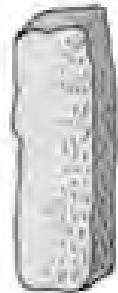
10x



5x



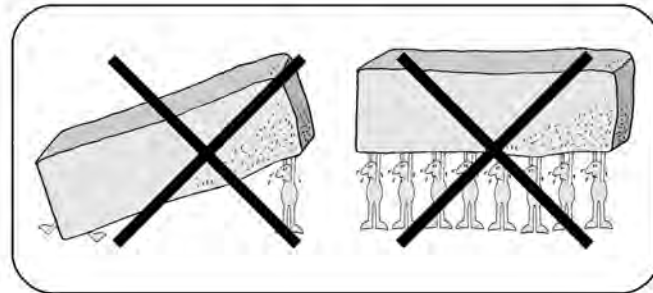
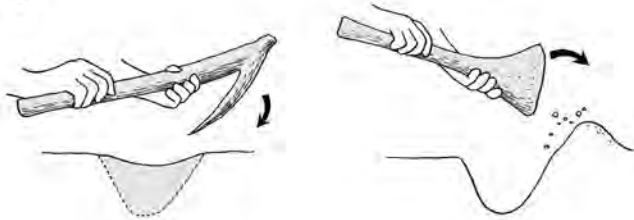
1x



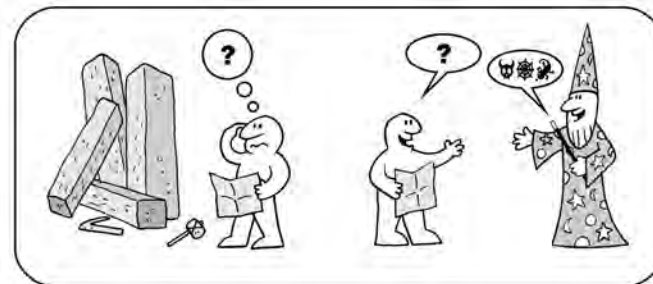
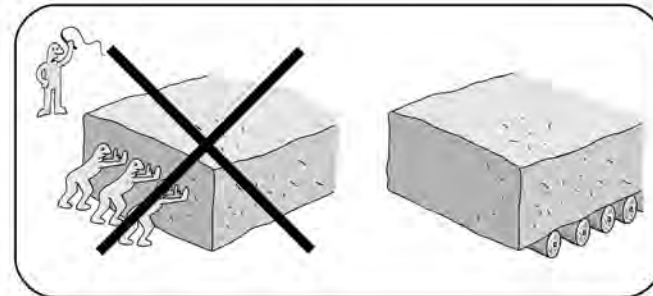
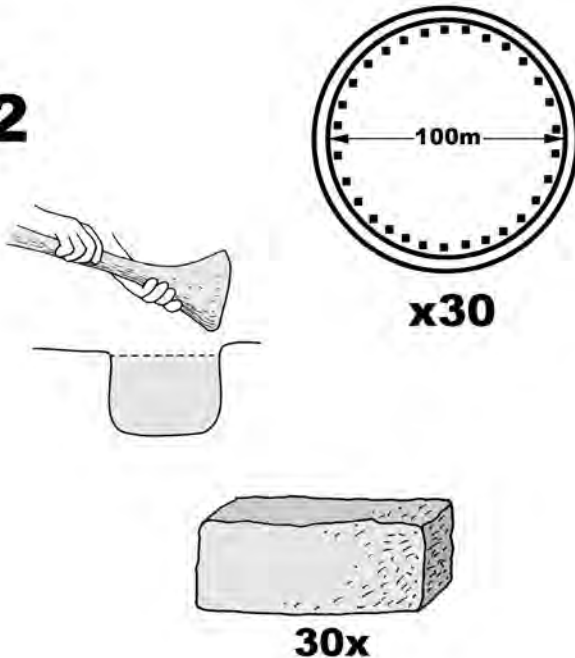
3x

Parts interact → they need to be assembled to work

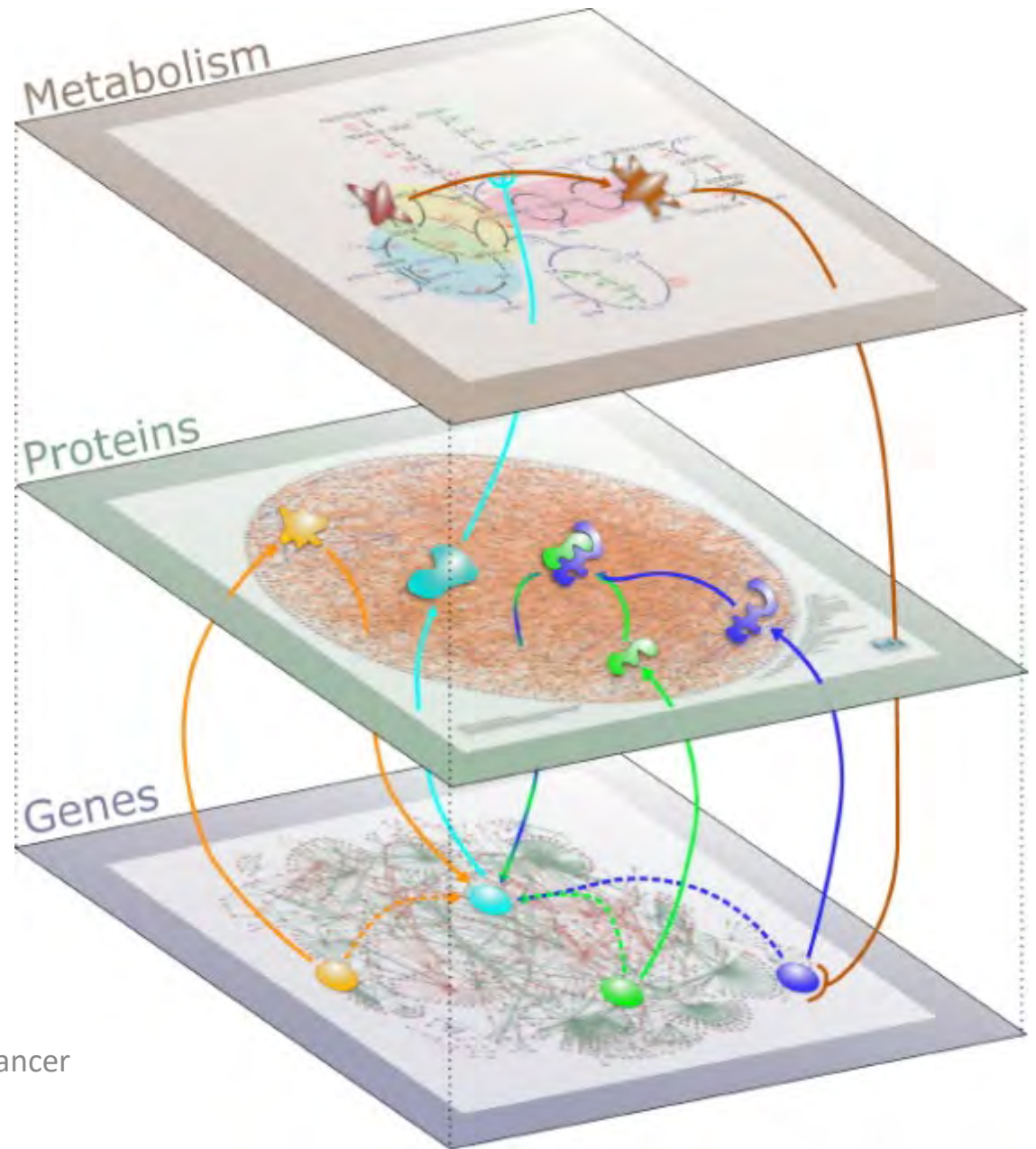
1



2



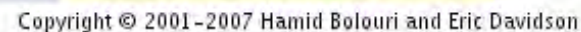
Intra-cellular Networks operate on multiple levels



Slides by Amitabh Sharma, PhD

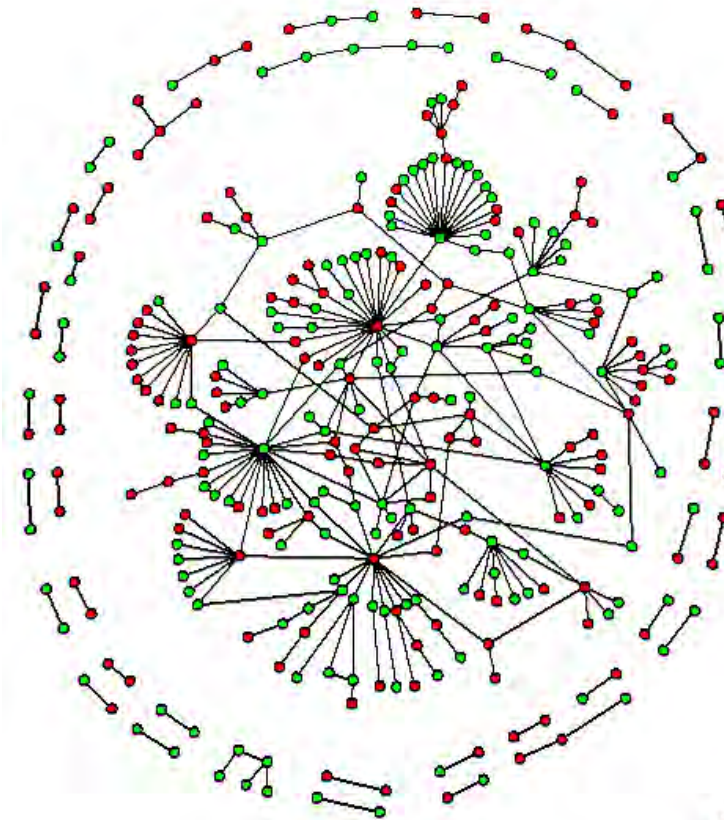
Northeastern University & Dana Farber Cancer
Institute

May 29, 2007

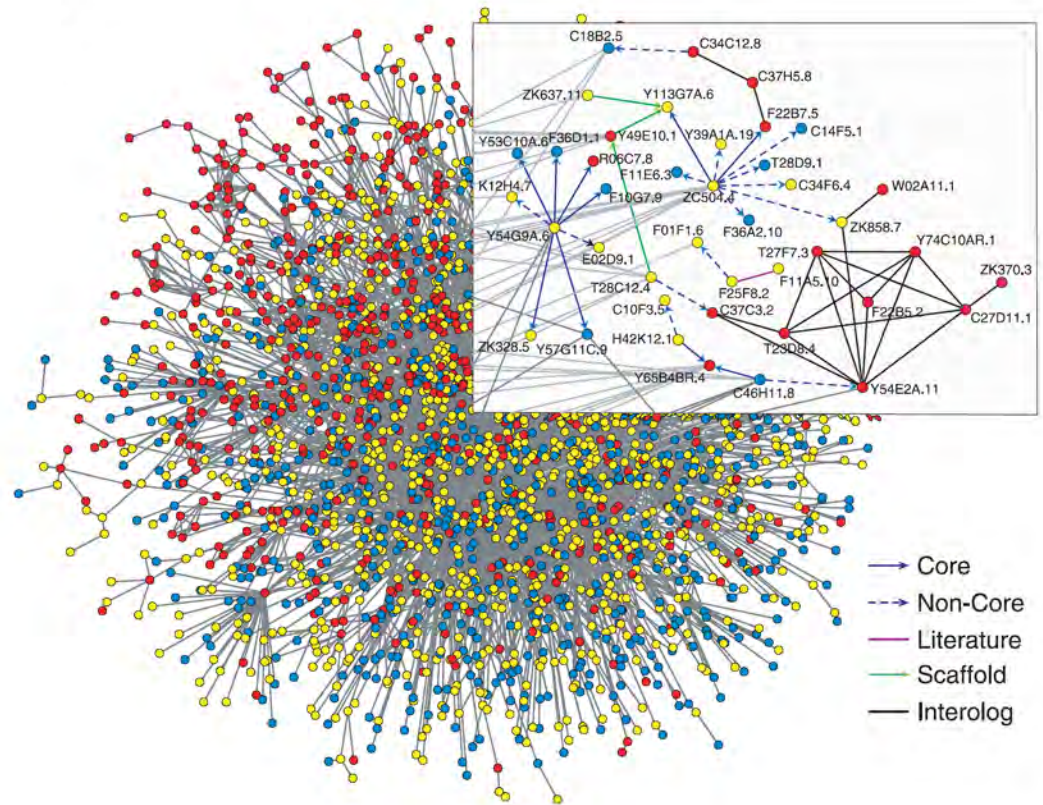


Ubiqu=ubiquitous; Mat = maternal; activ = activator; rep = repressor; unkn = unknown; Nucl. = nuclearization; χ = β -catenin source; n β -TCF = nuclearized b- β -catenin-Tcf1; ES = early signal; ECNS = early cytoplasmic nuclearization system; Zyg. N. = zygotic Notch

Interactions: 577,297 Proteins: 89,716

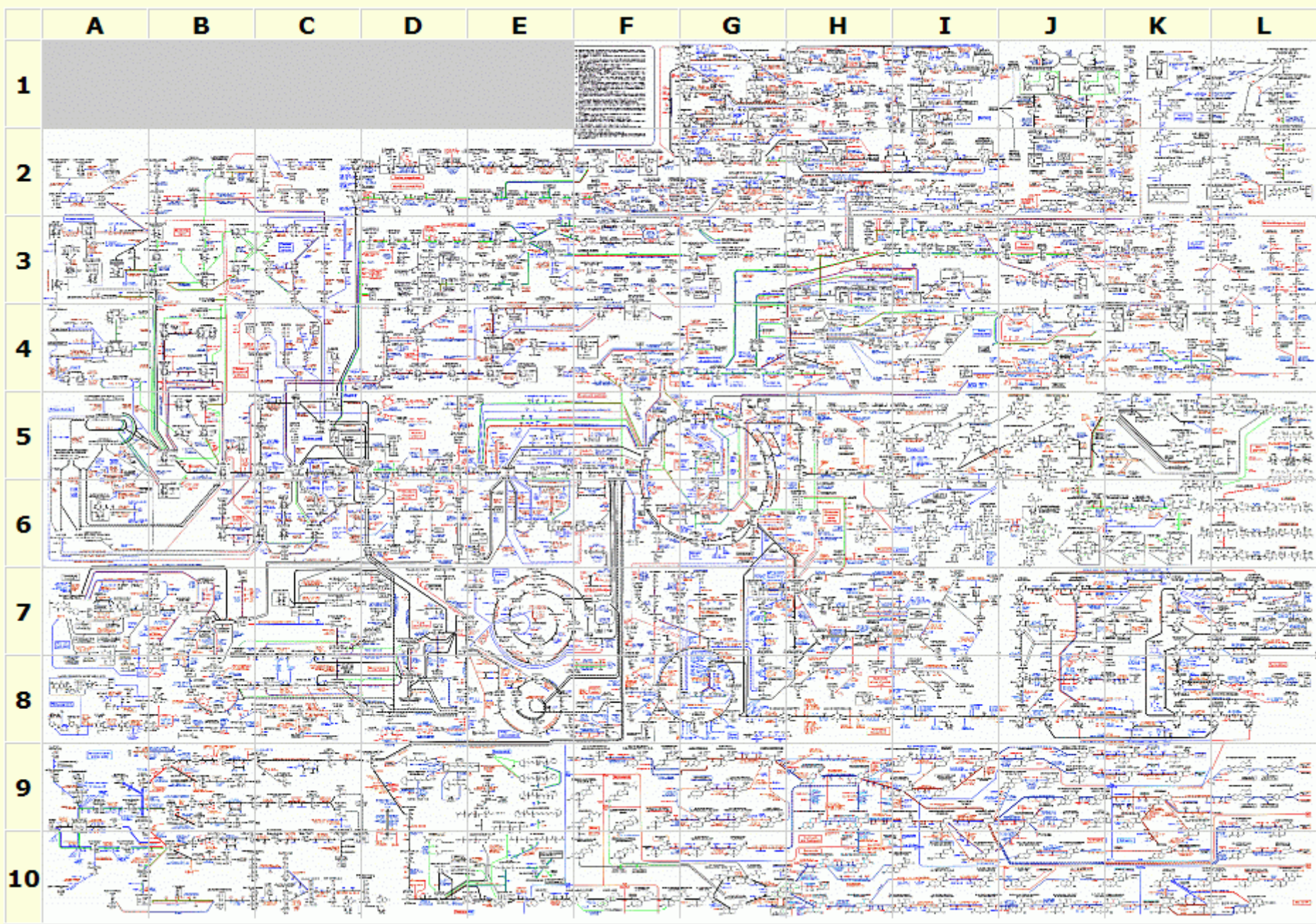


Baker's yeast *S. cerevisiae* (only nuclear proteins shown)
From S. Maslov, K. Sneppen, Science 2002



Worm *C. elegans*
From S. Lee et al , Science 2004

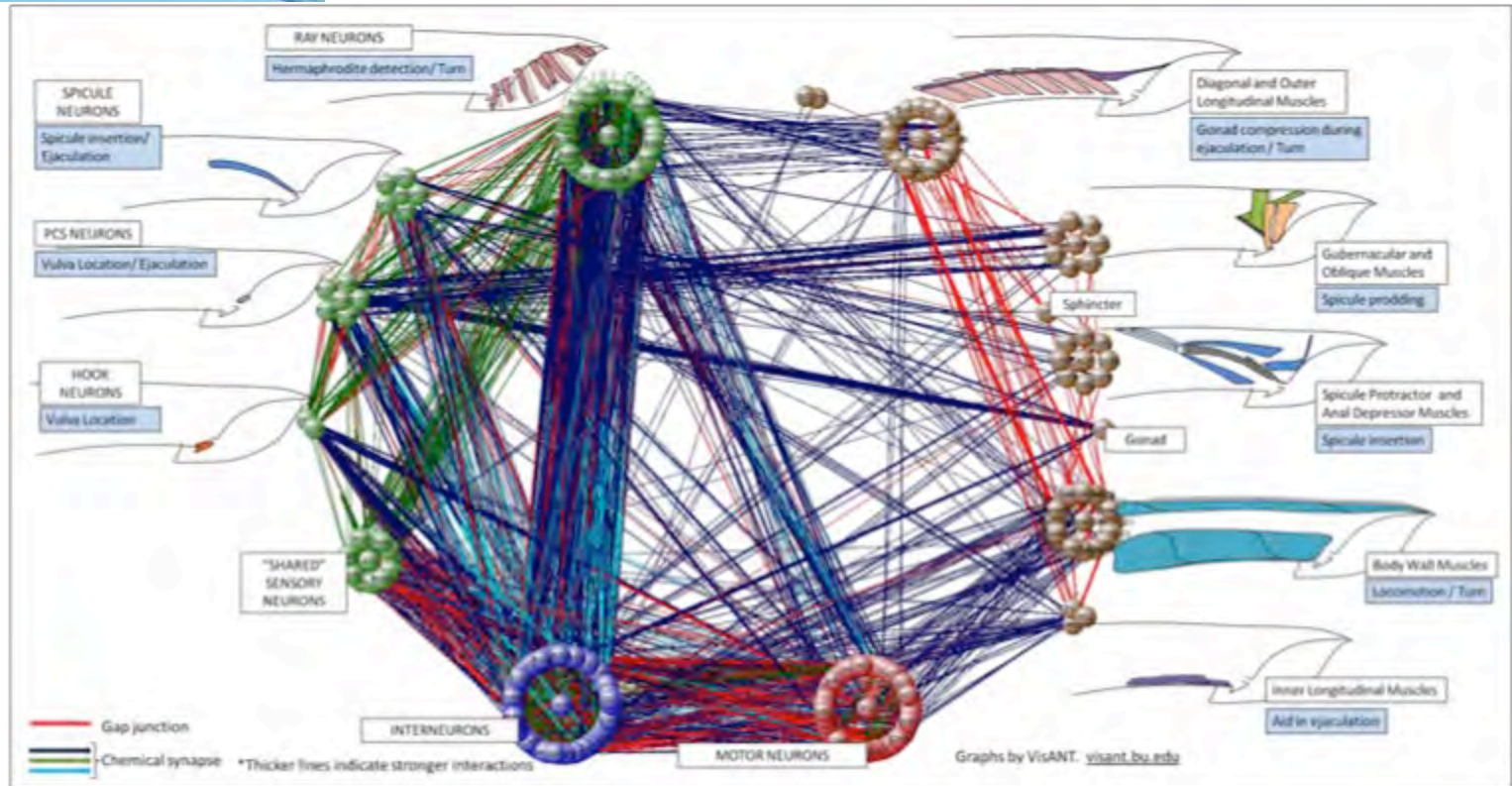
Metabolic pathway chart by ExPASy: 5702 reactions as of December 2015

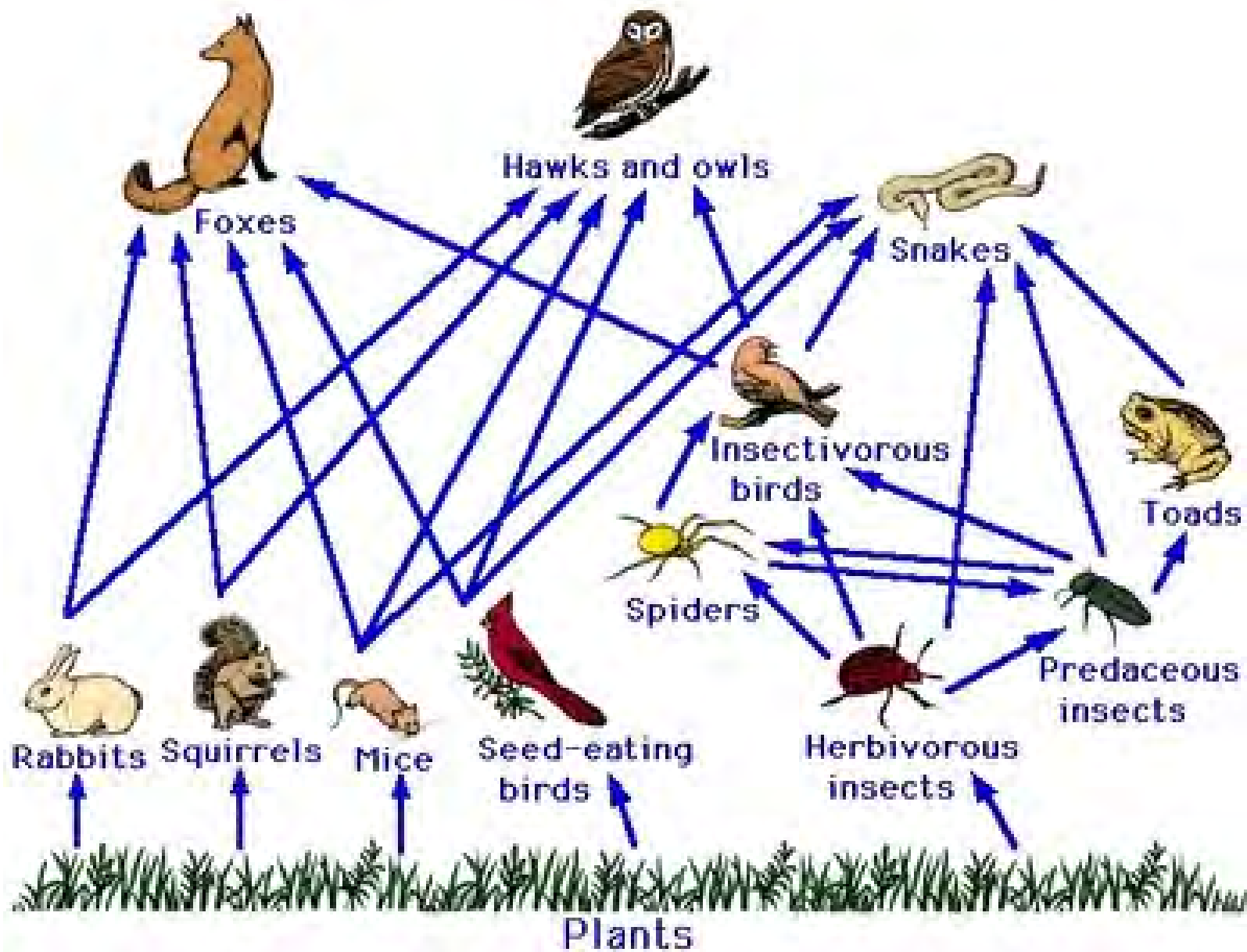


Brain and nerves of a worm



- Worm (*C. elegans*) has 302 neurons
- Our brain has 100 billion (10^{11}) neurons

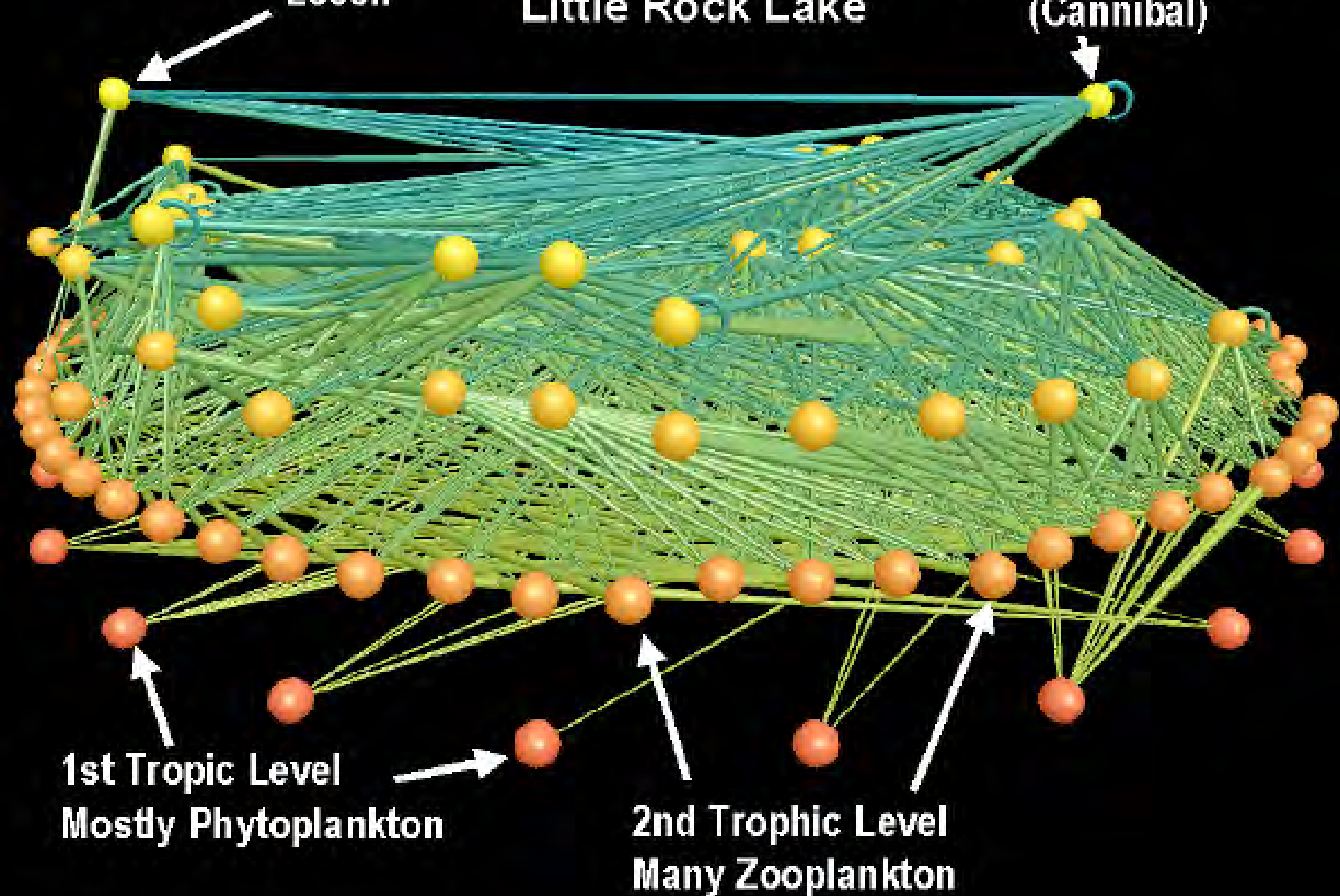




Food Web of Little Rock Lake

Smallmouth Bass
(Cannibal)

Leech



WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS
WHY DO I SAY UH

WHY IS SEA SALT BETTER
WHY ARE THERE TREES IN THE MIDDLE OF FIELDS
WHY IS THERE NOT A POKEMON MMO
WHY IS THERE LAUGHING IN TV SHOWS
WHY ARE THERE DOORS ON THE FREEWAY
WHY ARE THERE SO MANY SVCHOST.EXE RUNNING
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA
WHY ARE THERE SCARY SOUNDS IN MINECRAFT
WHY IS THERE KICKING IN MY STOMACH
WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES
WHY DO SNAKES EXIST
WHY DO OYSTERS HAVE PEARLS
WHY ARE DUCKS CALLED DUCKS
WHY DO THEY CALL IT THE CLAP
WHY ARE KYLE AND CARTMAN FRIENDS
WHY IS THERE AN ARROW ON AANG'S HEAD
WHY ARE TEXT MESSAGES BLUE
WHY ARE THERE MUSTACHES ON CLOTHES
WHY ARE THERE MUSTACHES ON CARS
WHY ARE THERE MUSTACHES EVERYWHERE
WHY ARE THERE SO MANY BIRDS IN OHIO
WHY IS THERE SO MUCH RAIN IN OHIO
WHY IS OHIO WEATHER SO WEIRD

WHY ARE THERE MALE AND FEMALE BIKES
WHY ARE THERE BRIDESMAIDS
WHY DO DYING PEOPLE REACH UP
WHY AREN'T THERE VARICOSE ARTERIES
WHY ARE OLD KLINGONS DIFFERENT



WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL

WHY DO TESTICLES MOVE
WHY ARE THERE PSYCHICS
WHY ARE HATS SO EXPENSIVE
WHY IS THERE CAFFEINE IN MY SHAMPOO
WHY DO YOUR BOOBS HURT

WHY DO IGUANAS DIE
WHY AREN'T THERE DINOSAUR GHOSTS

WHY AREN'T ECONOMISTS RICH
WHY DO AMERICANS CALL IT SOCCER
WHY ARE MY EARS RINGING
WHY ARE THERE SO MANY AVENGERS
WHY ARE THE AVENGERS FIGHTING THE X MEN
WHY IS WOLVERINE NOT IN THE AVENGERS

WHY IS EARTH TILTED
WHY IS SPACE BLACK
WHY IS OUTER SPACE SO COLD
WHY ARE THERE PYRAMIDS ON THE MOON
WHY IS NASA SHUTTING DOWN

WHY ARE THERE TINY SPIDERS IN MY HOUSE
WHY DO SPIDERS COME INSIDE
WHY ARE THERE HUGE SPIDERS IN MY HOUSE
WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE
WHY ARE THERE SPIDERS IN MY ROOM
WHY ARE THERE SO MANY SPIDERS IN MY ROOM
WHY DO SPIDER BITES ITCH
WHY IS DYING SO SCARY

WHY ARE THERE FEMALE MR NIMES

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY ARE THERE SLAVES IN THE BIBLE
WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY ARE THERE SWARMS OF GNATS
WHY IS THERE PHLEGM



WHY IS MT VESUVIUS THERE
WHY DO THEY SAY T MINUS
WHY ARE THERE OBELISKS
WHY ARE WRESTLERS ALWAYS WET
WHY ARE OCEANS BECOMING MORE ACIDIC
WHY IS ARWEN DYING
WHY AREN'T MY QUAIL LAYING EGGS
WHY AREN'T MY QUAIL EGGS HATCHING
WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT



WHY ARE THERE SO MANY CROWS IN ROCHESTER, MN
WHY IS PSYCHIC WEAK TO BUG
WHY DO CHILDREN GET CANCER
WHY IS POSEIDON ANGRY WITH ODYSSEUS
WHY IS THERE ICE IN SPACE

WHY ARE THERE DOGS AFRAID OF FIREWORKS
WHY IS THERE AN OWL IN MY BACKYARD
WHY IS THERE AN OWL OUTSIDE MY WINDOW
WHY IS THERE AN OWL ON THE DOLLAR BILL
WHY DO OWLS ATTACK PEOPLE
WHY ARE AK 47s SO EXPENSIVE
WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE
WHY ARE THERE GODS
WHY ARE THERE TWO SPOCKS
WHY ARE MY BOOBS ITCHY
WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG

Credit: XKCD
comics

Foundations of Probability

Random experiments

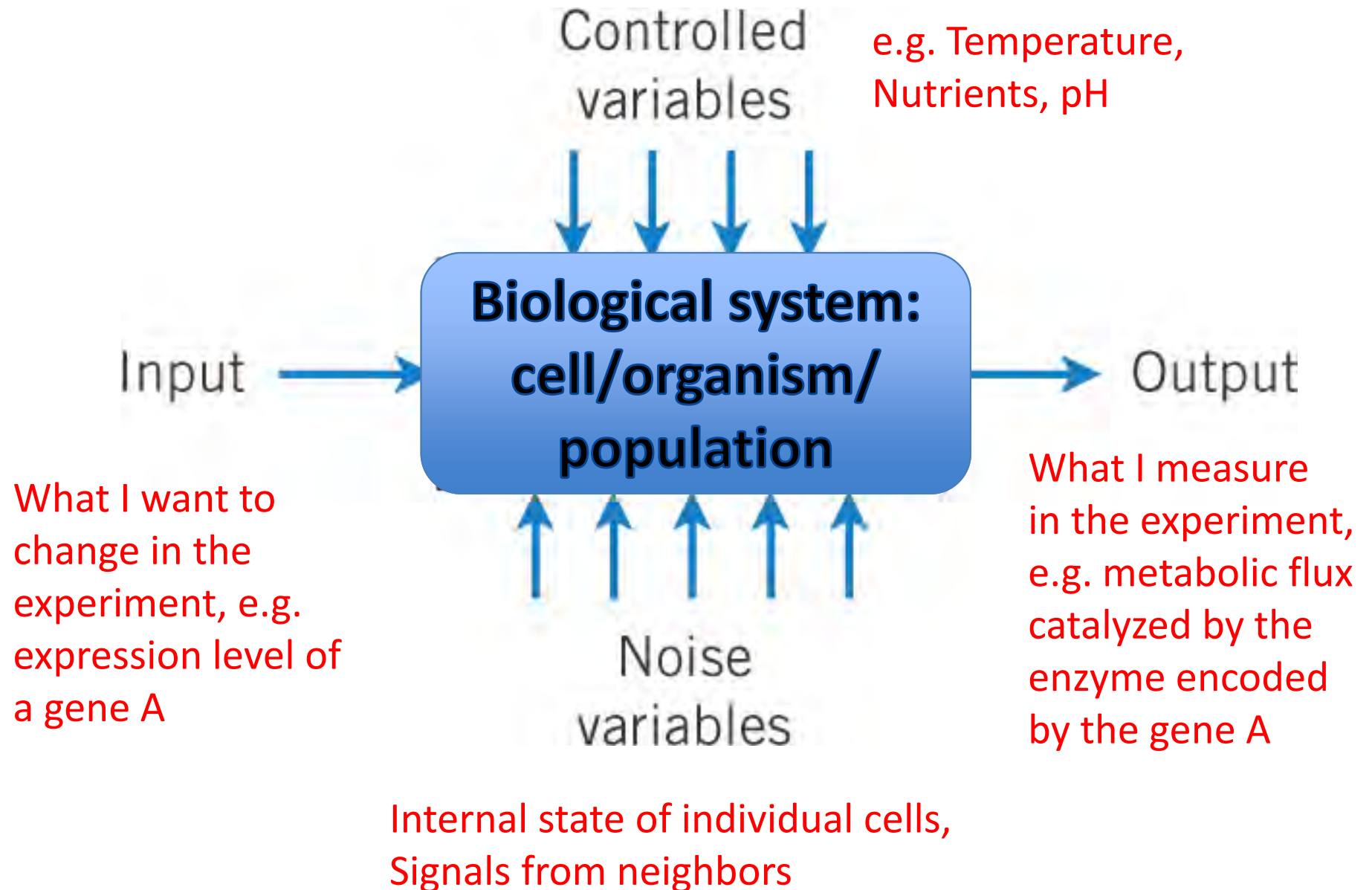
Sample spaces

Venn diagrams of
random events

Random Experiments

- An **experiment** is an operation or procedure, carried out under controlled conditions
 - Example: measure the metabolic flux through a reaction catalyzed by the enzyme A
- An experiment that can result in **different outcomes**, even if repeated in the same manner every time, is called a **random experiment**
 - Cell-to-cell variability due to history/genome variants
 - Noise in external parameters such as temperature, nutrients, pH, etc.
- **Evolution** offers ready-made random experiments
 - Genomes of different species
 - Genomes of different individuals within a species
 - Individual cancer cells

Variability/Noise Produce Output Variation



Sample Spaces

- Random experiments have unique outcomes.
- The set of all possible outcomes of a random experiment is called the sample space, S .
- S is discrete if it consists of a finite or countable infinite set of outcomes.
- S is continuous if it contains an interval (either a finite or infinite width) of real numbers.

Examples of a Sample Space

- Experiment measuring the abundance of mRNA expressed from a single gene

$S = \{x | x \geq 0\}$: continuous.

- Bin it into four groups

$S = \{\textit{below 10}, \textit{10-30}, \textit{30-100}, \textit{above 100}\}$: discrete.

- Is gene “on” (mRNA above 30)?

$S = \{\textit{true}, \textit{false}\}$: logical/Boolean/discrete.

Event

An event (E) is a **subset of the sample space** of a random experiment, i.e., **one or more** outcomes of the sample space.

- The **union** of two events is the event that consists of all outcomes that are contained in either of the two events. We denote the union as $E_1 \cup E_2$
- The **intersection** of two events is the event that consists of all outcomes that are contained in both of the two events. We denote the intersection as $E_1 \cap E_2$
- The **complement** of an event in a sample space is the set of outcomes in the sample space that are not in the event. We denote the complement of the event E as E' (sometimes E^c or \bar{E})

Examples

Discrete

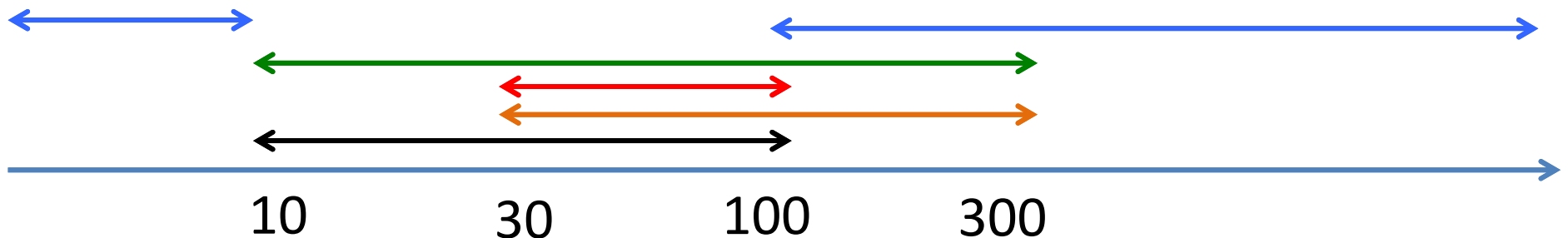
1. Assume you toss a coin once. The sample space is $S = \{H, T\}$, where H = head and T = tail and the event of a head is $\{H\}$.
2. Assume you toss a coin twice. The sample space is $S = \{(H, H), (H, T), (T, H), (T, T)\}$, and the event of obtaining exactly one head is $\{(H, T), (T, H)\}$.

Continuous

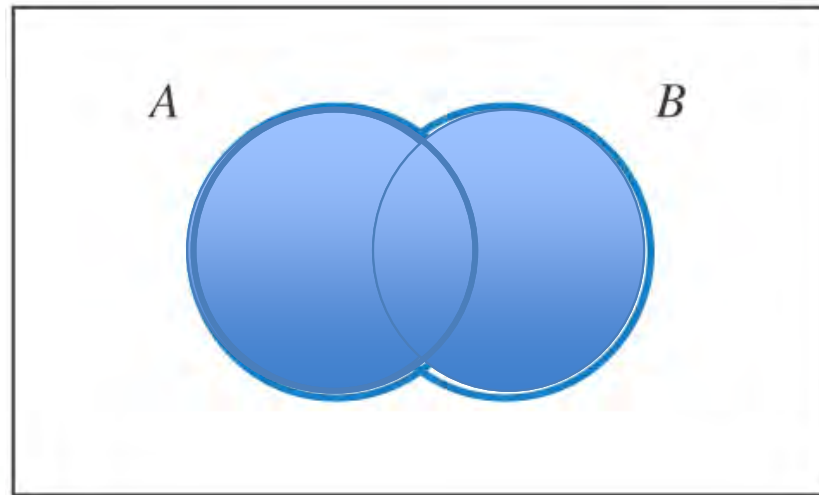
Sample space for the expression level of a gene: $S = \{x | x \geq 0\}$

Two events:

- $E1 = \{x | 10 < x < 100\}$
- $E2 = \{x | 30 < x < 300\}$
- $E1 \cap E2 = \{x | 30 < x < 100\}$
- $E1 \cup E2 = \{x | 10 < x < 300\}$
- $E1' = \{x | x \leq 10 \text{ or } x \geq 100\}$



Venn diagrams



$A \cup B$

S



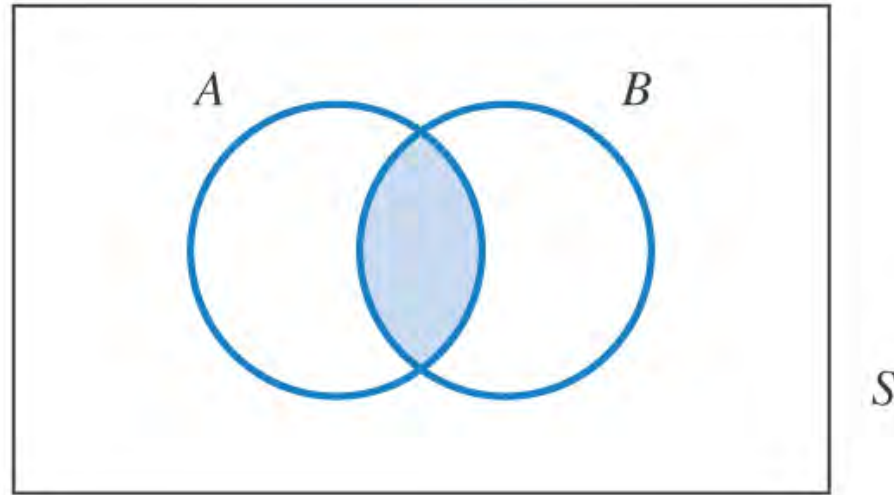
Find
5 differences
in beard and
hairstyle



John Venn (1843-1923)
British logician

John Venn (1990-)
Brooklyn hipster

Venn diagrams

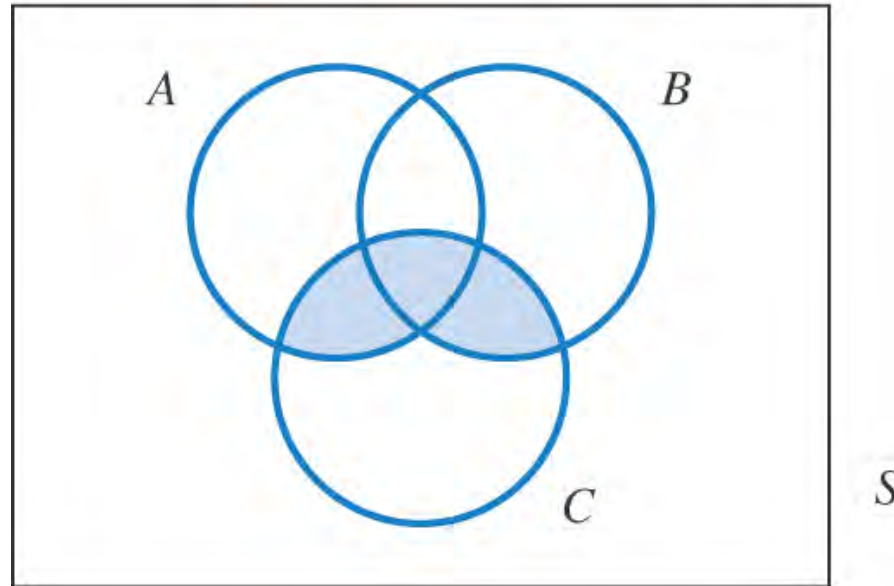


Which formula describes the blue region?

- A. $A \cup B$
- B. $A \cap B$
- C. A'
- D. B'

Get your i-clickers

Venn diagrams

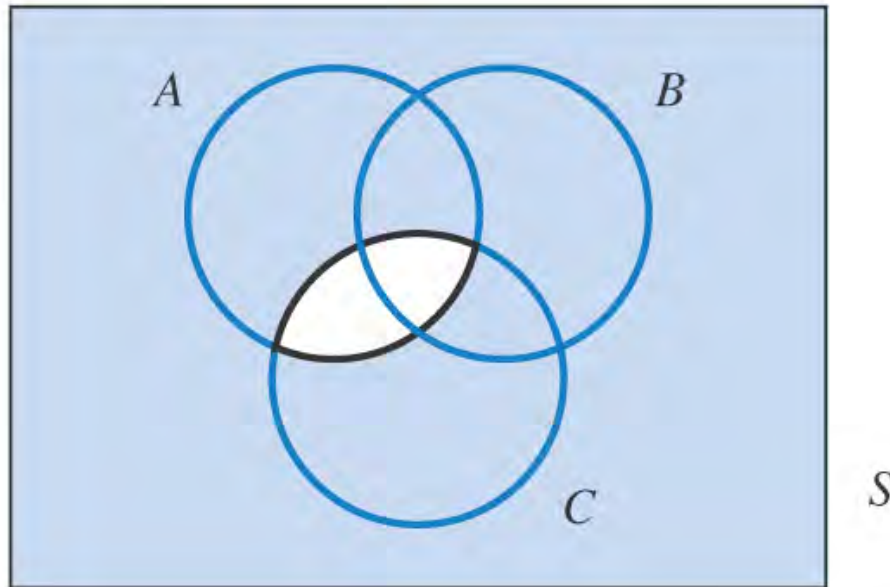


Which formula describes the blue region?

- A. $(A \cup B) \cap C$
- B. $(A \cap B) \cap C$
- C. $(A \cup B) \cup C$
- D. $(A \cap B) \cup C$

Get your i-clickers

Venn diagrams



Which formula describes the blue region?

- A. $A \cap C$
- B. $A' \cup C'$
- C. $(A \cap B \cap C)'$
- D. $(A \cap B) \cap C$

Get your i-clickers

Definitions of Probability

Two definitions of probability

- (1) **STATISTICAL PROBABILITY**: the relative frequency with which an event occurs in the long run
- (2) **INDUCTIVE PROBABILITY**: the degree of belief which it is reasonable to place in a proposition on given evidence

Statistical Probability

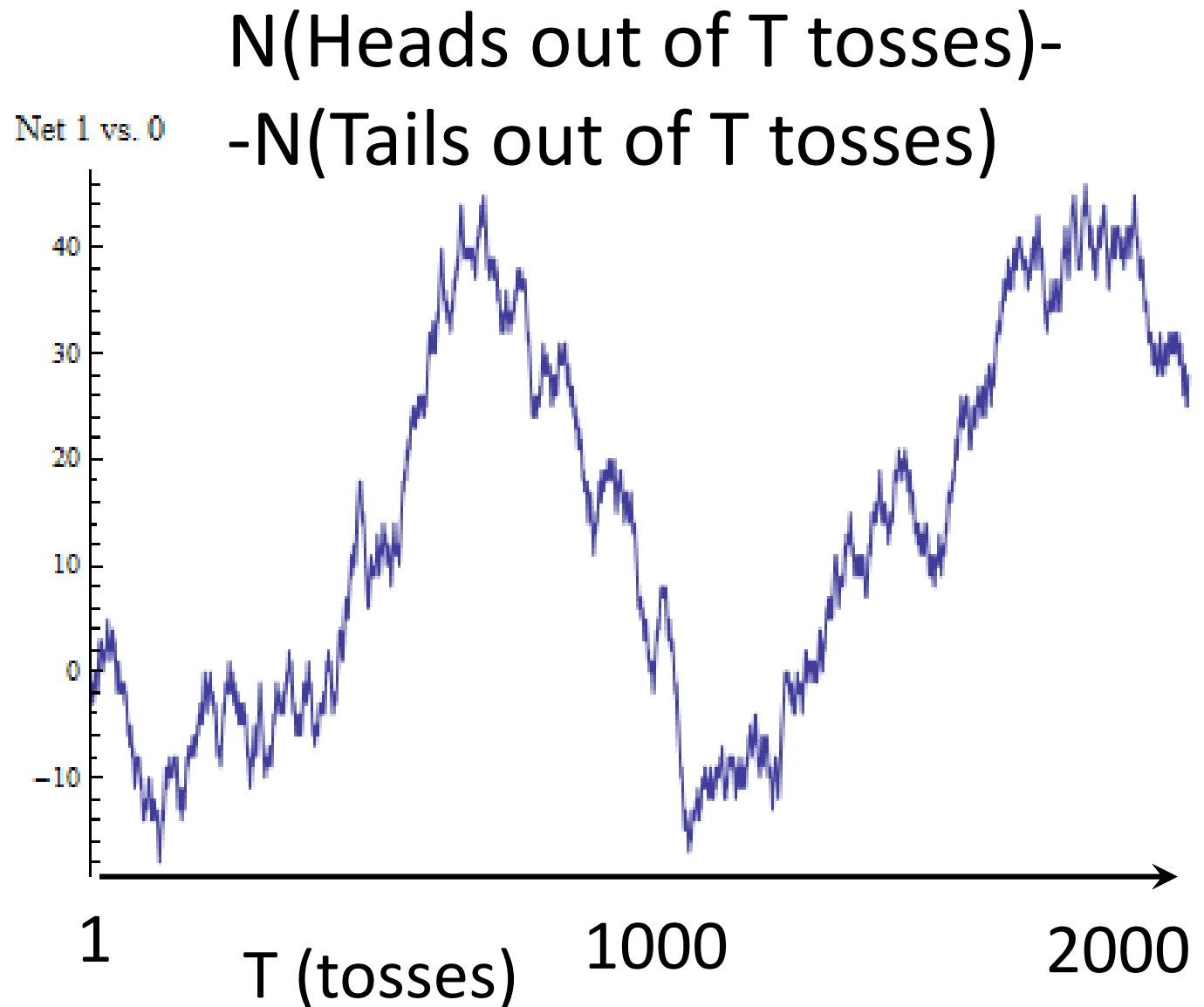
A **statistical probability** of an event is the **limiting value** of the **relative frequency** with it occurs in a **very large number** of **independent trials**

Empirical

Statistical Probability of a Coin Toss



John Edmund Kerrich
(1903–1985)
British/South African
mathematician

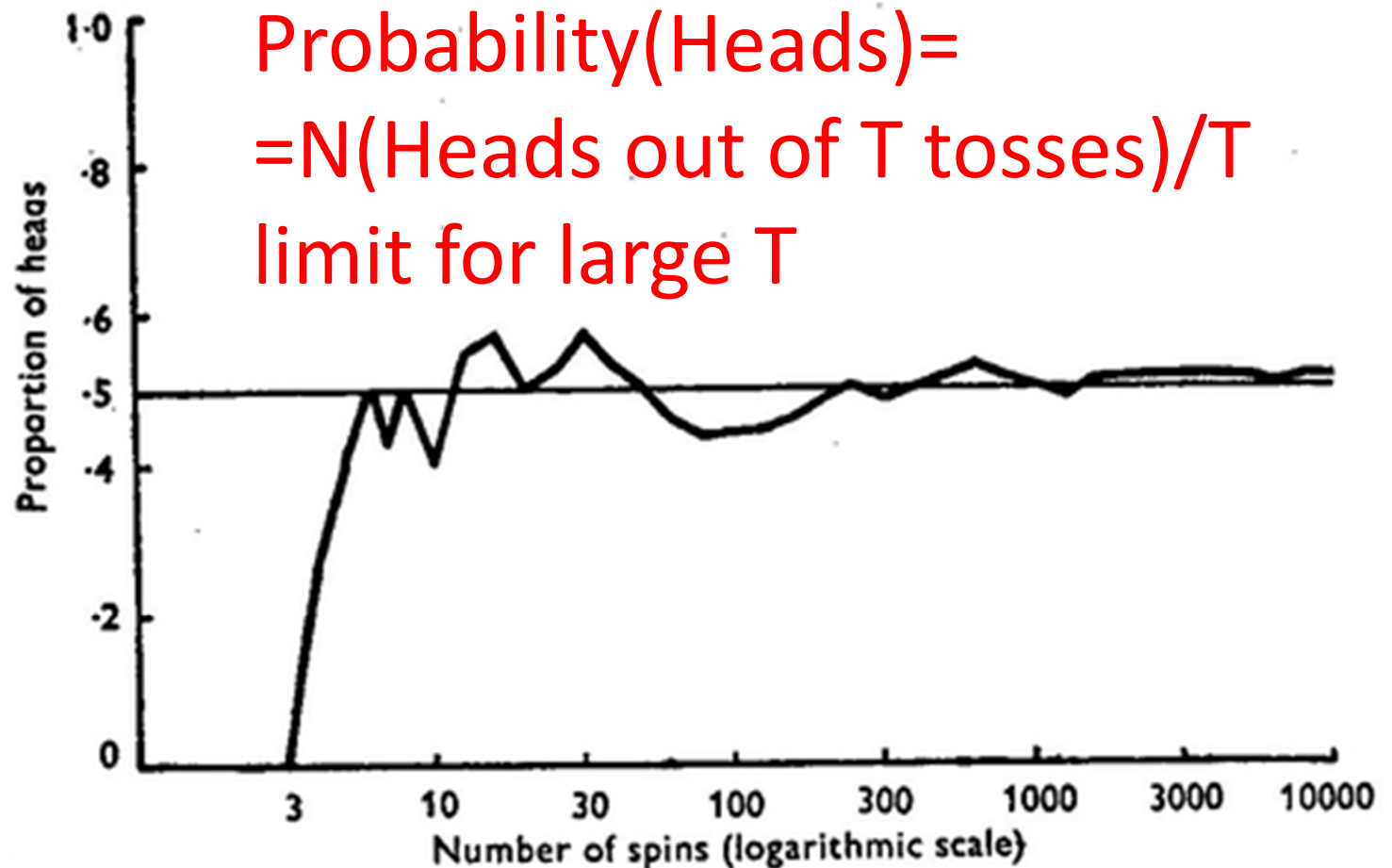


Excess of heads among 2,000 coin tosses (Kerrich 1946)

Statistical Probability of a Coin Toss



John Edmund Kerrich
(1903–1985)
British/South African
mathematician



Proportion of heads among 10,000 coin tosses (Kerrich 1946)

Who is ready to use Matlab?

- A. I have Matlab installed on my laptop
- B. I am ready to use Matlab on EWS
- C. I don't have it ready but plan to install it
- D. I am not ready but plan to use EWS
- E. I plan to use other software (Python, R, etc.)

Get your i-clickers

Matlab is easy to learn

- Matlab is the lingua franca of all of **engineering**
- Use online tutorials e.g.:
<https://www.youtube.com/watch?v=82TGgQApFIQ>
- **Matlab** is designed to work with **Matrices** → symbols ***** and **/** are understood as **matrix multiplication** and **division**
- Use **.*** and **./** for regular (non-matrix) multiplication
- Add **;** in the end of the line to avoid displaying the output on the screen
- **Loops**: **for** i=1:100; f(i)=floor(2.*rand); end;
- **Conditional statements**: **if** rand>0.5; count=count+1; end;
- **Plotting**: **plot**(x,y,'ko-'); or **semilogx**(x,y,'ko-'); or **loglog**(x,y,'ko-'); .
To keep **adding plots onto the same axes** use: **hold on**;
To **create a new axes** use **figure**;
- **Generating matrices**: **rand**(100) – generates square matrix 100x100.
Confusing! Use **rand**(100,1) or **zeros**(30,20), or **randn**(1,40) (Gaussian);
- If Matlab complains multiplying matrices **check sizes** using **whos** and if needed **use transpose** operation: **x=x'**;

A Matlab Cheat-sheet (MIT 18.06, Fall 2007)

Basics:

save 'file.mat'	save variables to <i>file.mat</i>
load 'file.mat'	load variables from <i>file.mat</i>
diary on	record input/output to file <i>diary</i>
diary off	stop recording
whos	list all variables currently defined
clear	delete/undefine all variables
help command	quick help on a given <i>command</i>
doc command	extensive help on a given <i>command</i>

Defining/changing variables:

x = 3	define variable <i>x</i> to be 3
x = [1 2 3]	set <i>x</i> to the 1×3 row-vector (1,2,3)
x = [1 2 3];	same, but don't echo <i>x</i> to output
x = [1;2;3]	set <i>x</i> to the 3×1 column-vector (1,2,3)
A = [1 2 3 4; 5 6 7 8; 9 10 11 12];	
	set <i>A</i> to the 3×4 matrix with rows 1,2,3,4 etc.
x(2) = 7	change <i>x</i> from (1,2,3) to (1,7,3)
A(2,1) = 0	change <i>A</i> _{2,1} from 5 to 0

Arithmetic and functions of numbers:

3*4, 7+4, 2-6 8/3	multiply, add, subtract, and divide numbers
3^7, 3^(8+2i)	compute 3 to the 7th power, or 3 to the 8+2i power
sqrt(-5)	compute the square root of -5
exp(12)	compute e^{12}
log(3), log10(100)	compute the natural log (ln) and base-10 log (log ₁₀)
abs(-5)	compute the absolute value -5
sin(5*pi/3)	compute the sine of 5π/3
besselj(2,6)	compute the Bessel function $J_2(6)$

Arithmetic and functions of vectors and matrices:

x * 3	multiply every element of <i>x</i> by 3
x + 2	add 2 to every element of <i>x</i>
x + y	element-wise addition of two vectors <i>x</i> and <i>y</i>
A * y	product of a matrix <i>A</i> and a vector <i>y</i>
A * B	product of two matrices <i>A</i> and <i>B</i>
x * y	not allowed if <i>x</i> and <i>y</i> are two column vectors!
x .* y	element-wise product of vectors <i>x</i> and <i>y</i>
A^3	the square matrix <i>A</i> to the 3rd power
x^3	not allowed if <i>x</i> is not a square matrix!
x.^3	every element of <i>x</i> is taken to the 3rd power
cos(x)	the cosine of every element of <i>x</i>
abs(A)	the absolute value of every element of <i>A</i>
exp(A)	e to the power of every element of <i>A</i>
sqrt(A)	the square root of every element of <i>A</i>
expm(A)	the matrix exponential e^A
sqrtm(A)	the matrix whose square is <i>A</i>

Transposes and dot products:

x.', A.'	the transposes of <i>x</i> and <i>A</i>
x', A'	the complex-conjugate of the transposes of <i>x</i> and <i>A</i>
x' * y	the dot (inner) product of two column vectors <i>x</i> and <i>y</i>
dot(x,y), sum(x.*y)	...two other ways to write the dot product
x * y'	the outer product of two column vectors <i>x</i> and <i>y</i>

Constructing a few simple matrices:

rand(12,4)	a 12×4 matrix with uniform random numbers in [0,1)
randn(12,4)	a 12×4 matrix with Gaussian random (center 0, variance 1)
zeros(12,4)	a 12×4 matrix of zeros
ones(12,4)	a 12×4 matrix of ones
eye(5)	a 5×5 identity matrix I ("eye")
eye(12,4)	a 12×4 matrix whose first 4 rows are the 4×4 identity
linspace(1.2,4.7,100)	row vector of 100 equally-spaced numbers from 1.2 to 4.7
7:15	row vector of 7,8,9,...,14,15
diag(x)	matrix whose diagonal is the entries of <i>x</i> (and other elements = 0)

Portions of matrices and vectors:

x(2:12)	the 2nd to the 12th elements of <i>x</i>
x(2:end)	the 2nd to the last elements of <i>x</i>
x(1:3:end)	every third element of <i>x</i> , from 1st to the last
x(:)	all the elements of <i>x</i>
A(5,:)	the row vector of every element in the 5th row of <i>A</i>
A(5,1:3)	the row vector of the first 3 elements in the 5th row of <i>A</i>
A(:,2)	the column vector of every element in the 2nd column of <i>A</i>
diag(A)	column vector of the diagonal elements of <i>A</i>

Solving linear equations:

A \ b	for <i>A</i> a matrix and <i>b</i> a column vector, the solution <i>x</i> to $Ax=b$
inv(A)	the inverse matrix A^{-1}
[L,U,P] = lu(A)	the LU factorization $PA=LU$
eig(A)	the eigenvalues of <i>A</i>
[V,D] = eig(A)	the columns of <i>V</i> are the eigenvectors of <i>A</i> , and the diagonals diag(<i>D</i>) are the eigenvalues of <i>A</i>

Plotting:

plot(y)	plot <i>y</i> as the <i>y</i> axis, with 1,2,3,... as the <i>x</i> axis
plot(x,y)	plot <i>y</i> versus <i>x</i> (must have same length)
plot(x,A)	plot columns of <i>A</i> versus <i>x</i> (must have same # rows)
loglog(x,y)	plot <i>y</i> versus <i>x</i> on a log-log scale
semilogx(x,y)	plot <i>y</i> versus <i>x</i> with <i>x</i> on a log scale
semilogy(x,y)	plot <i>y</i> versus <i>x</i> with <i>y</i> on a log scale
fplot(@(x) ...expression..., [a,b])	plot some expression in <i>x</i> from <i>x</i> = <i>a</i> to <i>x</i> = <i>b</i>
axis equal	force the <i>x</i> and <i>y</i> axes of the current plot to be scaled equally
title('A Title')	add a title <i>A Title</i> at the top of the plot
xlabel('blah')	label the <i>x</i> axis as <i>blah</i>
ylabel('blah')	label the <i>y</i> axis as <i>blah</i>
legend('foo','bar')	label 2 curves in the plot <i>foo</i> and <i>bar</i>
grid	include a grid in the plot
figure	open up a new figure window

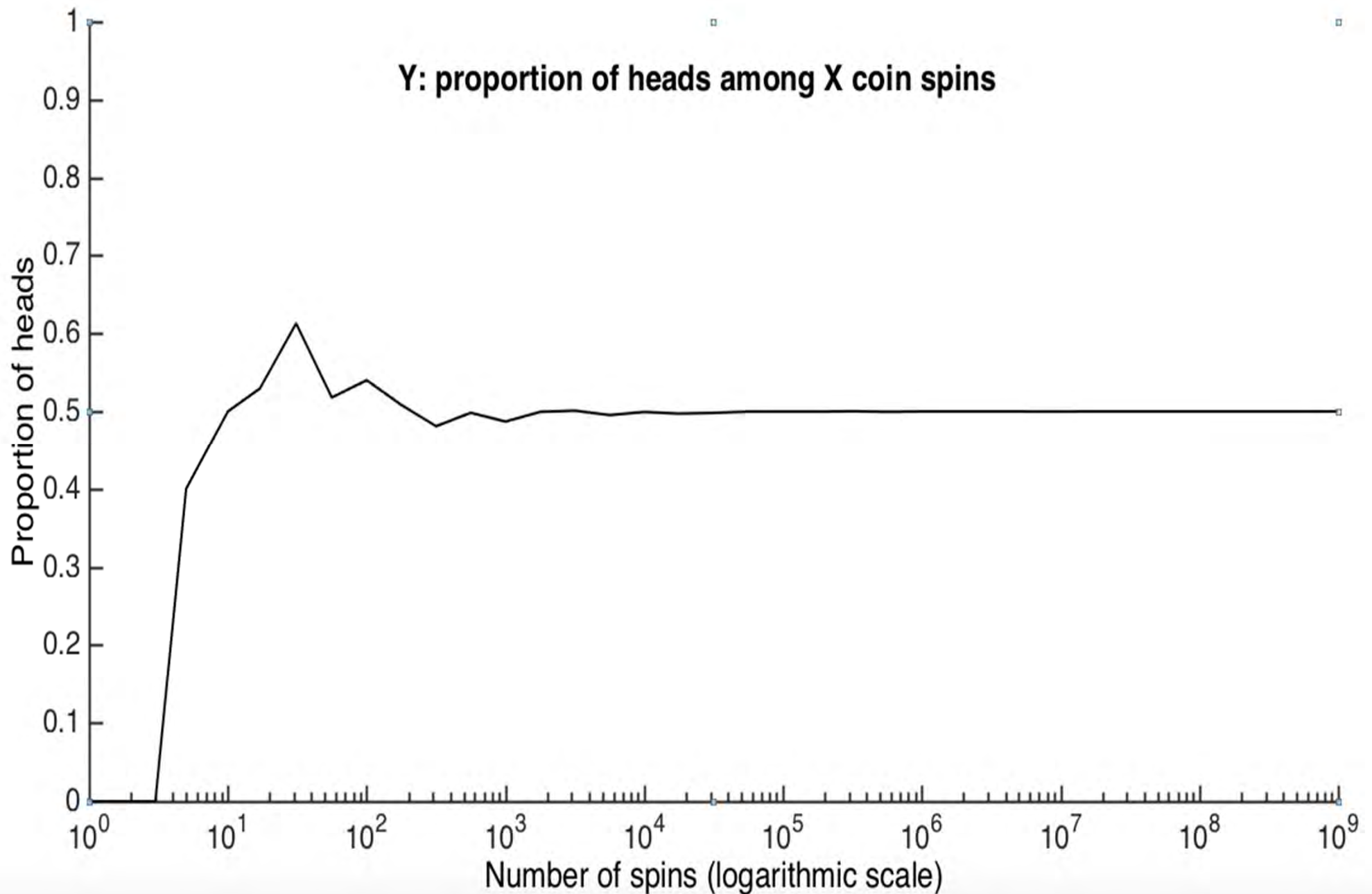
Matlab group exercise

Each table to edit the file `coin_toss_template.m` (replace all ?? with commands/variables/operations) or writes a new Matlab (Python, R, or anything else) script to:

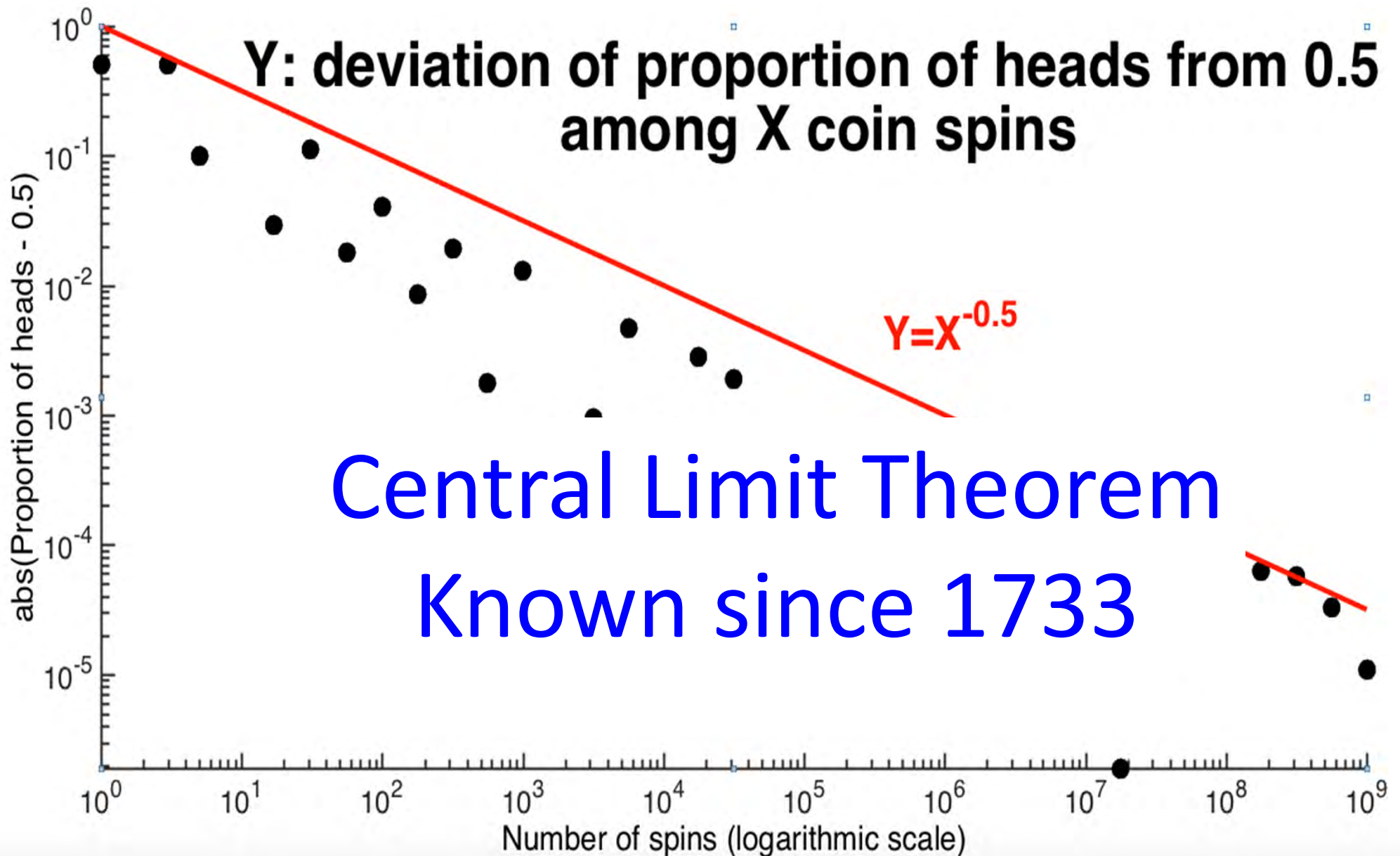
- Simulate a fair coin toss experiment
- Generate multiple tosses of a fair coin:
1 – heads, 0 - tails
- Calculate the fraction of heads ($f_heads(t)$) at timepoints:
t=10; 100; 1000; 10,000; 100,000; 1,000,000;10,000,000
coin tosses
- Plot fraction of heads $f_heads(t)$ vs t with a logarithmic t-axis
- Plot $abs(f_heads(t)-0.5)$ vs t on a log-log plot (both axes are logarithmic)

How I did it

- Stats=1e7;
- r0=rand(Stats,1); r1=floor(2.*r0);
- n_heads(1)=r1(1);
- for t=2:Stats; n_heads(t)=n_heads(t-1)+r1(t); end;
- tp=[1, 10,100,1000, 10000, 100000, 1000000, 10000000]
- np=n_heads(tp); fp=np./tp
- figure; semilogx(tp,fp,'ko-');
- hold on; semilogx([1,10000000],[0.5,0.5],'r--');
- figure; loglog(tp,abs(fp-0.5),'ko-');
- hold on; loglog(tp,0.5./sqrt(tp),'r--');



Proportion of heads among 1,000,000,000 coin tosses
(10^5 more than Kerrich) took me 33 seconds on my Surface Book

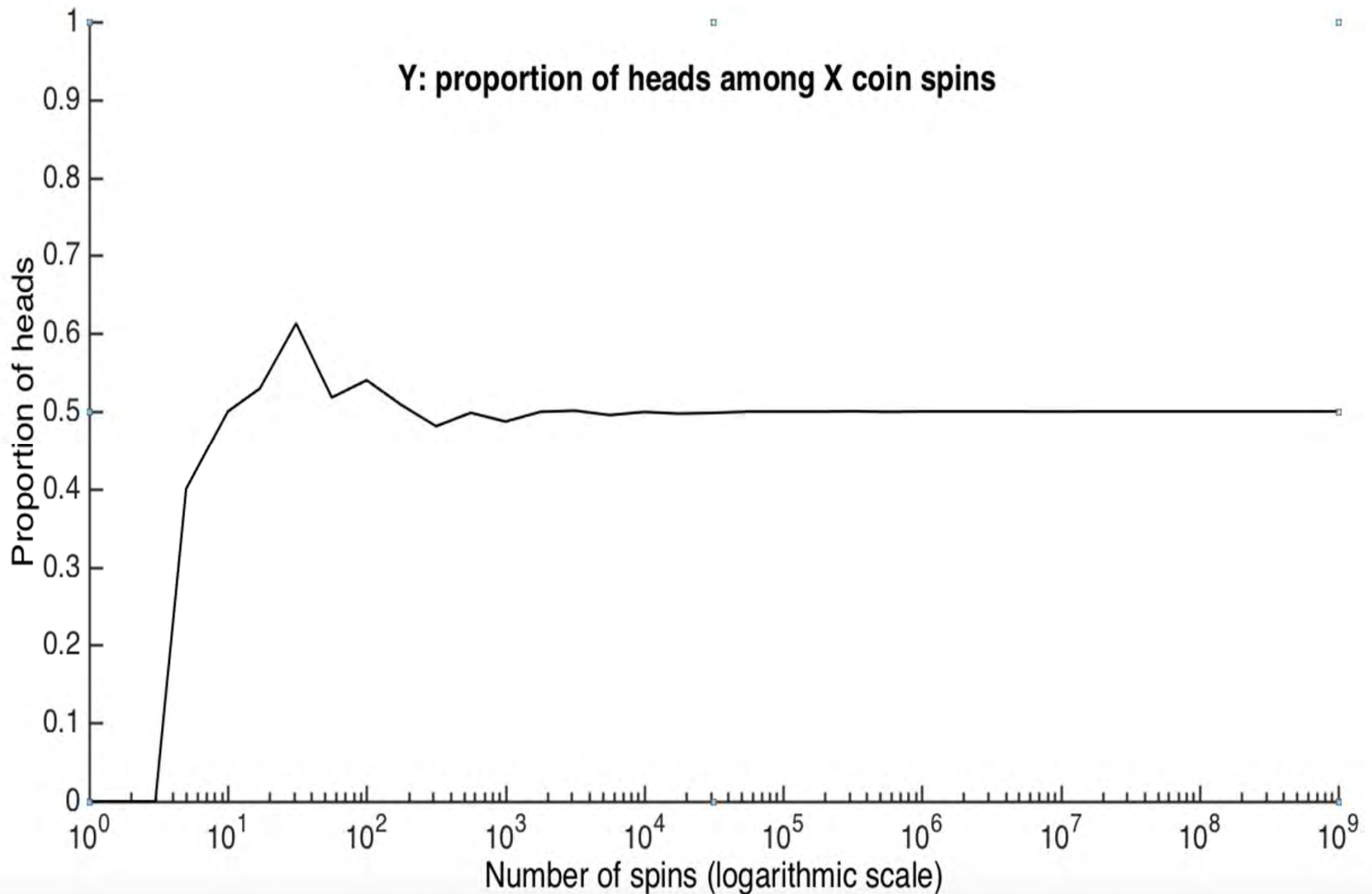


ABS(Proportion of heads-0.5)
among 100,000,000 coin tosses

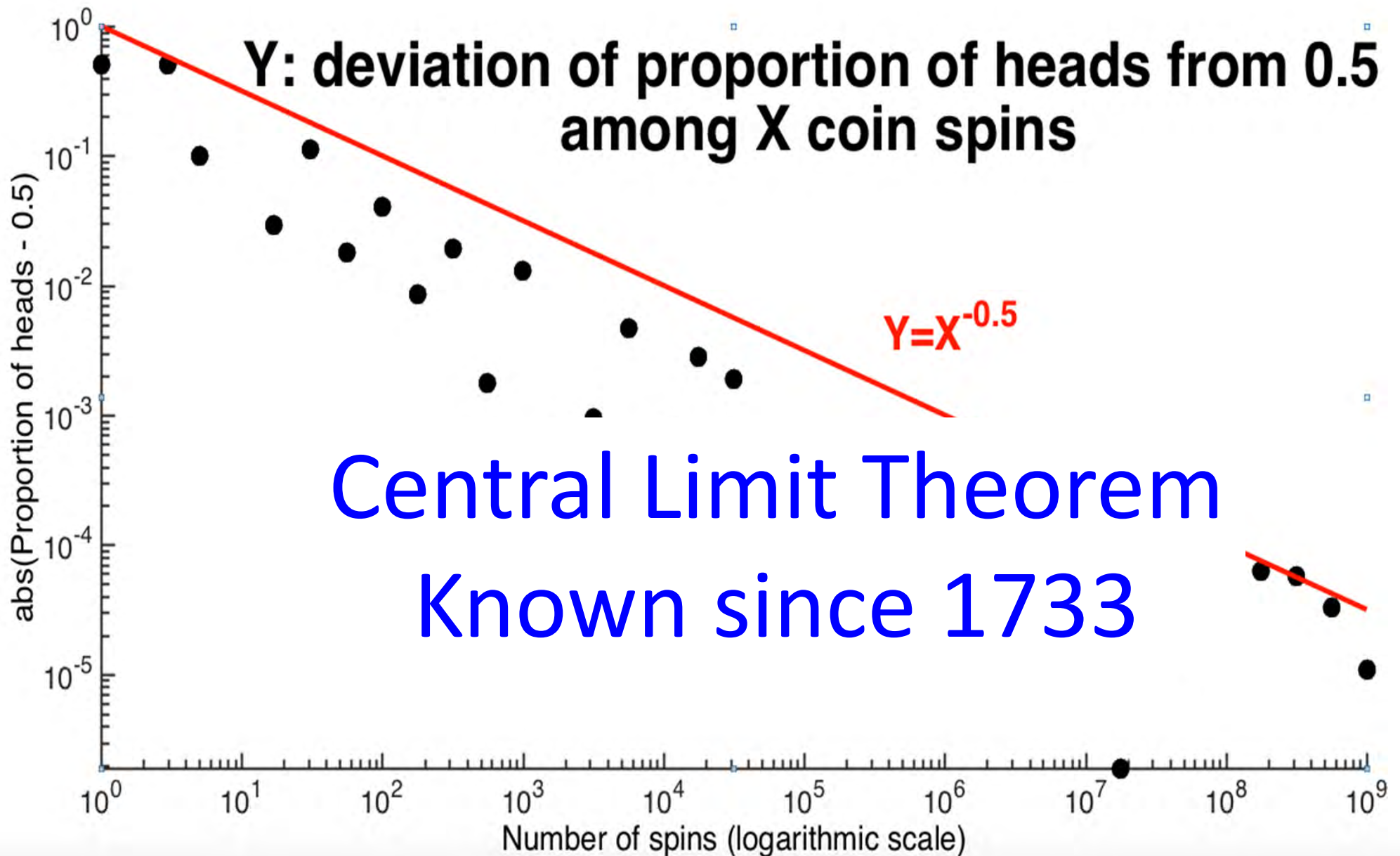
Matlab group exercise

Each table to edit the file `coin_toss_template.m` (replace all ?? with commands/variables/operations) or writes a new Matlab (Python, R, or anything else) script to:

- Simulate a fair coin toss experiment
- Generate multiple tosses of a fair coin:
1 – heads, 0 - tails
- Calculate the fraction of heads ($f_heads(t)$) at timepoints:
t=10; 100; 1000; 10,000; 100,000; 1,000,000;10,000,000
coin tosses
- Plot fraction of heads $f_heads(t)$ vs t with a logarithmic t-axis
- Plot $abs(f_heads(t)-0.5)$ vs t on a log-log plot (both axes are logarithmic)



Proportion of heads among 1,000,000,000 coin tosses
(10^5 more than Kerrich) took me 33 seconds on my Surface Book



ABS(Proportion of heads-0.5)
among 100,000,000 coin tosses

Definitions of Probability

Two definitions of probability

- (1) **STATISTICAL PROBABILITY**: the relative frequency with which an event occurs in the long run
- (2) **INDUCTIVE PROBABILITY**: the degree of belief which it is reasonable to place in a proposition on given evidence

Inductive Probability

An inductive probability of an event the degree of belief which it is rational to place in a hypothesis or proposition on given evidence.

Logical

Principle of indifference

- **Principle of Indifference** states that two **events are equally probable** if we have **no reason to suppose** that one of them will happen rather than the other. (Laplace, 1814)

- Unbiased coin:
probability Heads =
probability Tails = $\frac{1}{2}$

- Symmetric die:
probability of each side = $\frac{1}{6}$

**Pierre-Simon,
marquis de Laplace**
(1749 –1827)
French mathematician,
physicist, astronomer



Inductive = Naïve probability

- If space S is finite and **all outcomes are equally likely**, then

$$\text{Prob}(\text{Event } E) = \frac{\# \text{ of outcomes in } E}{\# \text{ of all outcomes in } S}$$

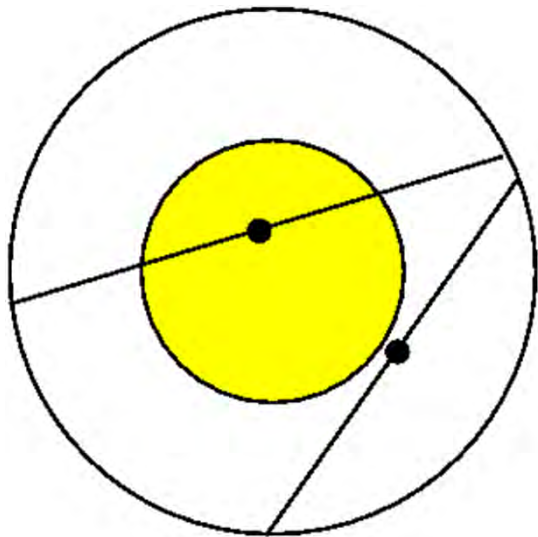
- Can also work with continuous is $\#$ is replaced with Area or Volume
- Unbiased coin: $\text{Prob}(\text{Heads}) = \text{Prob}(\text{Tails}) = 1/2$
- Symmetric die: probability of each side = $1/6$
- Lottery outcomes are not symmetric: It is not a 50%-50% chance to win or loose in a lottery

Inductive probability can lead to trouble

- Glass contains a mixture of wine and water and proportion of water to wine can be anywhere between 1:1 and 2:1
- (i) We can argue that the proportion of water to wine is equally likely to lie between 1 and 1.5 as between 1.5 and 2.
- (ii) Consider now ratio of wine to water. It is between 0.5 and 1. Based on the same argument it is equally likely in $[1/2, 3/4]$ as it is in $[3/4, 1]$. But then water to wine ratio is equally likely to lie between 1 and $4/3=1.333...$ as it is to lie between 1.333.. and 2. This is clearly inconsistent with the previous calculation...
- Paradox solved by clearly defining the experimental design:
 - For (i) use fixed amount of wine (1 liter) and select a uniformly-distributed random number between 1 and 2 for water.
 - For (ii) use 1 liter of water and select uniformly-distributed a random number between 0.5 and 1 for wine.
 - Different experiments – different answers
- Paradox is old. It is attributed to (among others) Joseph Bertrand

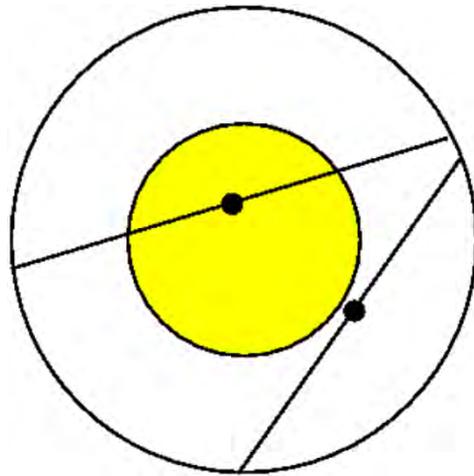
Better known Bertrand's paradox

- Take a circle of radius 2 and randomly draw a line segment through the circle. What is the probability P that the line intersects a concentric circle of radius 1?



Joseph Bertrand
(1822 –1900)
French mathematician

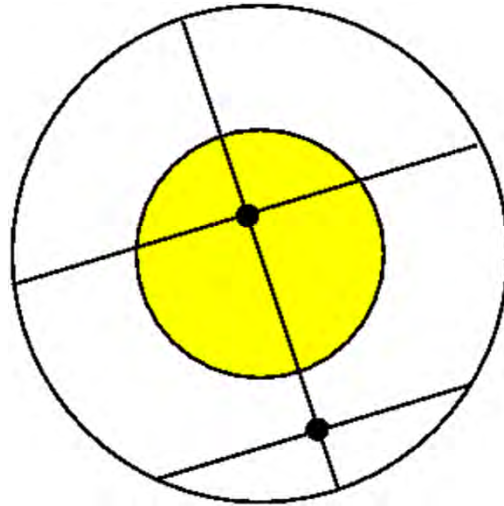
Solution #1



1. **Random point in 2D:** Each line has a unique midpoint, and a line will intersect the inner circle if its midpoint lies inside inner circle. Thus, P = probability that a randomly chosen midpoint lies in the inner circle:

$$P = \frac{\text{Area of the inner circle}}{\text{Area of the outer circle}} = \frac{\pi}{\pi 2^2} = \frac{1}{4}.$$

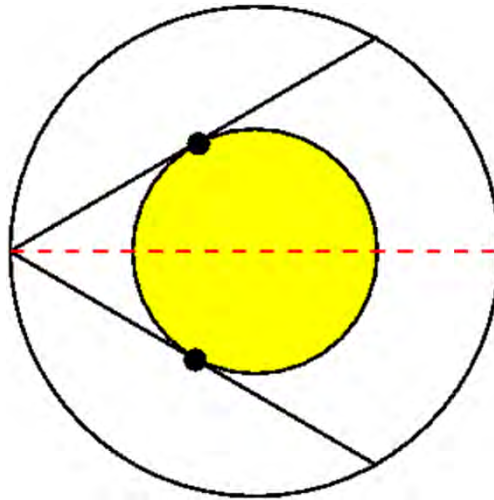
Solution #2



2. **Random point along the diameter:** Each line has a unique perpendicular bisector of length 4. So, P = probability that the midpoint lies on the inner part of the diameter:

$$P = \frac{\text{Length of the inner part of the diameter}}{\text{Length of the diameter}} = \frac{2}{4} = \frac{1}{2}$$

Solution #3

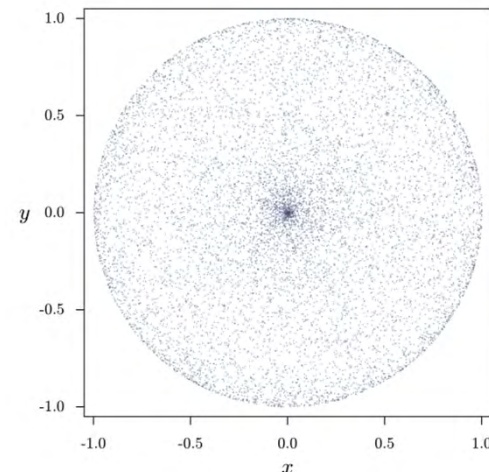
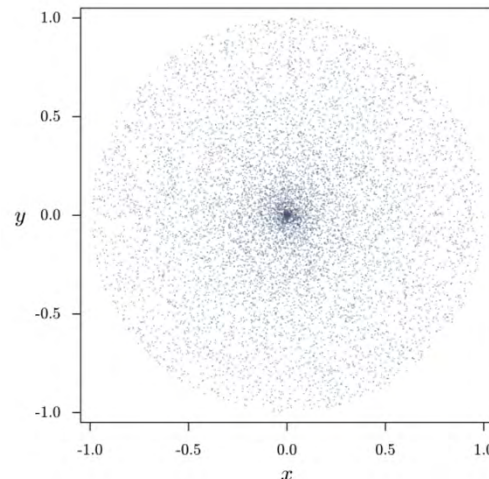
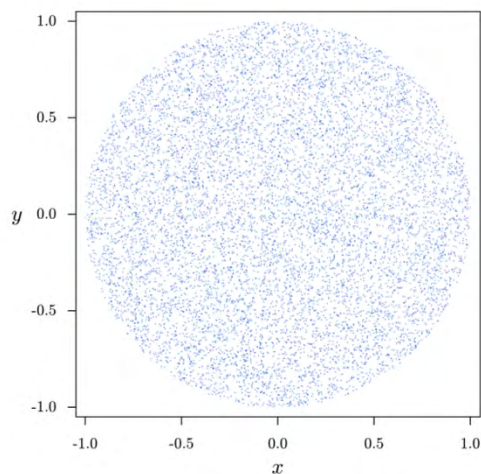


3. **Random angle:** Whether a line intersects the inner circle is determined by the angle it makes with the diameter intersecting the line on the outer circle:

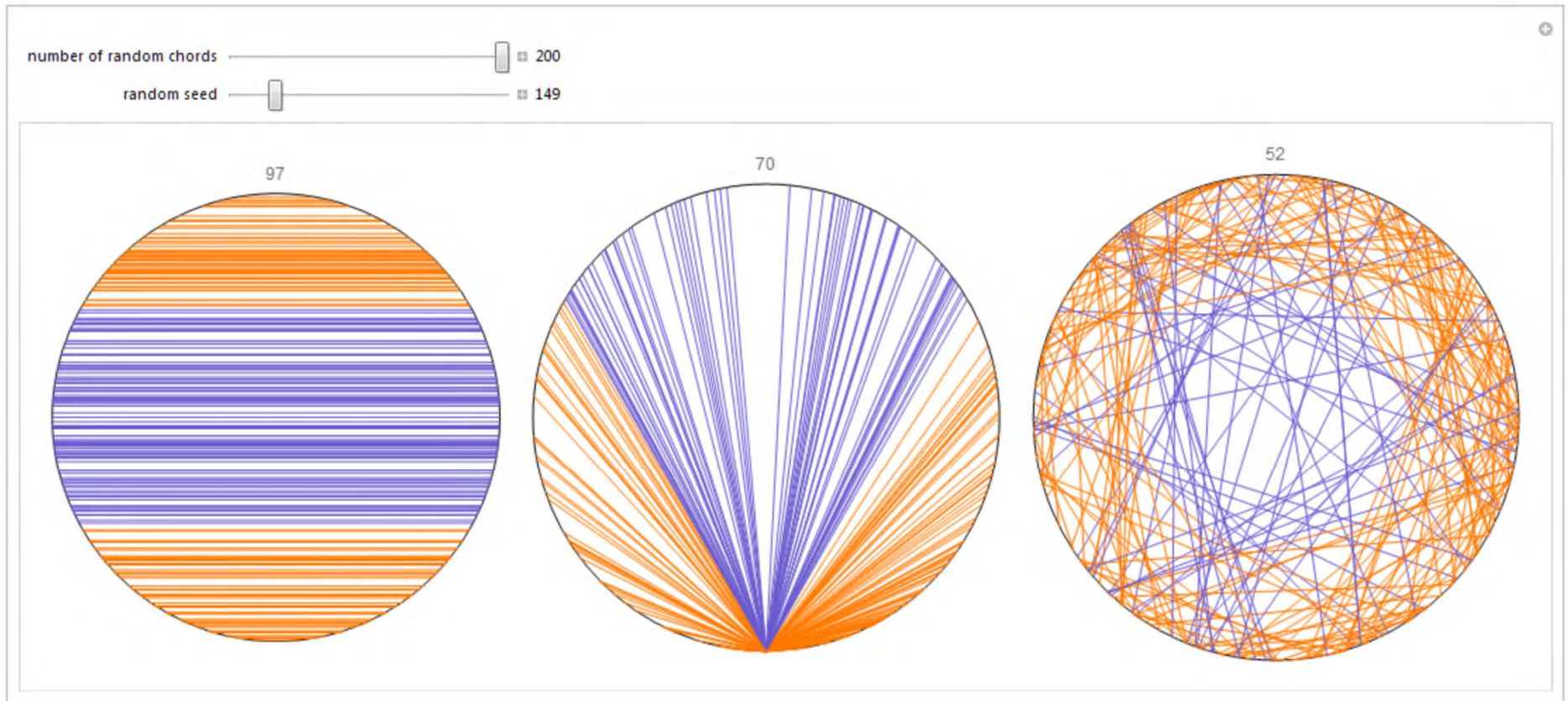
$$P = \frac{\pi/6}{\pi/2} = \boxed{\frac{1}{3}}.$$

So, is probability $1/4$, $1/2$, or $1/3$?

- Depends on how a “random” arc is selected:
 - For #1: select a point inside big circle and then draw an arc with this point as the center. Prob= $1/4$
 - For #2: select a diameter and a point on this diameter, then draw an arc. Prob= $1/2$
 - For #3: select a point on the circle and random angle. Prob= $1/3$



Mathematica visualization



I have two children.

One of them is a boy born on Tuesday.

What is the probability I have two boys?

A. $1/2$

B. $1/3$

C. $2/3$

D. $13/27$

E. I don't know

Get your i-clickers

Inductive probability
relies on combinatorics
or the art of counting
combinations

Counting – Multiplication Rule

- Multiplication rule:

- Let an operation consist of k steps and

- n_1 ways of completing the step 1,
- n_2 ways of completing the step 2, ... and

.....

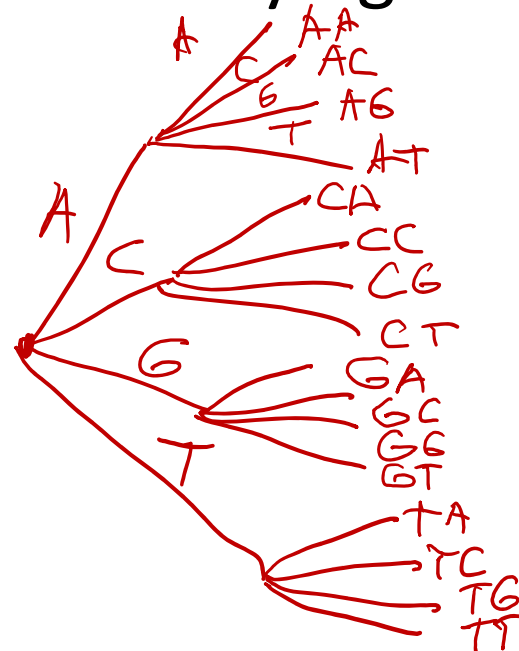
- n_k ways of completing the step k .

- Then, the total number of ways of carrying the entire operation is:

- $n_1 * n_2 * ... * n_k$

$$n_1 = n_2 = 4$$

Example: DNA 2-mer



- $S = \{A, C, G, T\}$ the set of 4 DNA bases
 - Number of k-mers is $4^k = 4 * 4 * 4 \dots * 4$ (k –times)
 - There are $4^3 = 64$ triplets in the genetic code
 - There are only 20 amino acids (AA)+1 stop codon
 - There is redundancy: same AA coded by 1-3 codons
 - Evidence of natural selection: “silent” changes of bases are more common than AA changing ones
- A protein-coding part of the gene is typically 1000 bases long
 - There are $4^{1000} = 2^{2000} \sim 10^{600}$ possible sequences of **just one gene**
 - Or $(10^{600})^{25,000} = 10^{15,000,000}$ of 25,000 human genes.
 - For comparison, the Universe has between 10^{78} and 10^{80} atoms and is $4 * 10^{17}$ seconds old.

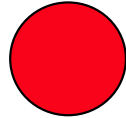
Counting – Permutation Rule

- A permutation is a unique sequence of distinct items.
- If $S = \{a, b, c\}$, then there are 6 permutations
 - Namely: abc, acb, bac, bca, cab, cba (**order matters**)
- # of permutations for a set of n items is $n!$
- $n!$ (factorial function) = $n * (n-1) * (n-2) * \dots * 2 * 1$
- $7! = 7 * 6 * 5 * 4 * 3 * 2 * 1 = 5,040$
- By definition: $0! = 1$

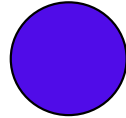
Multiplication and permutation
rules are two examples
of a general
problem, where
a sample of size k is drawn
from a population of
 n distinct objects

Balls drawn from an urn (or bowl)

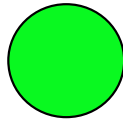
1 ball is red



1 ball is blue



1 ball is green

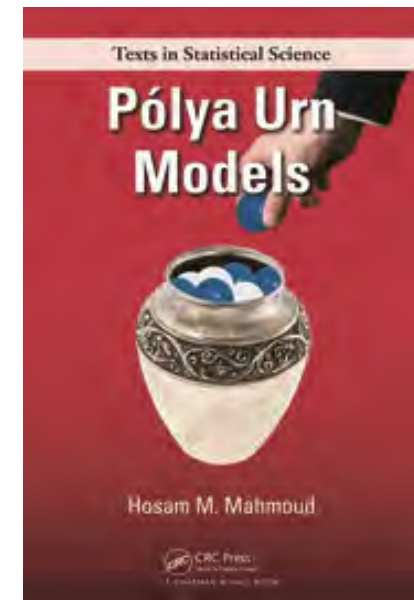


$n=3$ balls of different colors in an urn from which I draw $k=2$ balls one at a time

- Do I put each ball back to the bag after drawing it?
 - Yes: problem with replacement
 - No: problem without replacement
- Do I keep track of the order in which balls are drawn?
 - Yes: the order matters
 - No: the order does not matter

George Pólya

- George Pólya (December 13, 1887 – September 7, 1985) was a Hungarian mathematician. He was a professor of mathematics from 1914 to 1940 at ETH Zürich and from 1940 to 1953 at Stanford University. He made fundamental contributions to combinatorics, number theory, numerical analysis and probability theory.

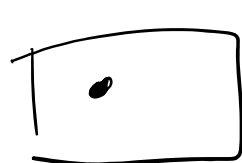


How many ways to Choose a sample of K objects out of a population of n objects

	Order matters	Order does not matter
replace	$n \times n \times n \times \dots \times n$ $= n^K$	$\frac{n^K}{K!}$ not all objects are different
Do not replace	$n \times (n-1) \times$ $\times (n-2) \times \dots \times$ $(n-K+1) =$ $= \frac{n!}{(n-K)!}$	All objects are different \rightarrow $\frac{n!}{(n-K)!} \times \frac{1}{K!} = \binom{n}{K}$

How to solve the problem of K out of n
with replacement but where order
does not matter?

Let's solve $n=2$ problem first:



object 1



object 2

$K=3$

4 possibilities



(1)



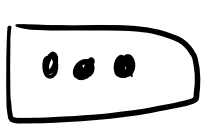
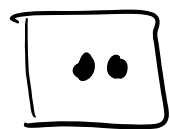
(2)



(3)



(4)



$n=4, K=7$



K dots, $n-1$ box boundaries

$$\binom{K+n-1}{K} = \frac{(K+n-1)!}{K! (n-1)!}$$

ways to distribute

Sampling table

How many ways to choose a sample of k objects out of population of n objects?

	Order matters	Order does not matter
Replacement	$(n)^k$	Difficult: $\binom{n+k-1}{k} = \frac{(n+k-1)!}{(n-1)!k!}$
No replacement	$n(n-1)(n-2)\dots(n-k+1) = \frac{n!}{(n-k)!}$	$\binom{n}{k} = \frac{n!}{(n-k)!k!}$

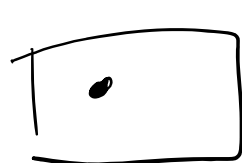
Inductive probability
relies on combinatorics
or the art of counting
combinations

How many ways to Choose a sample of K objects out of a population of n objects

	Order matters	Order does not matter
replace	$n \times n \times n \times \dots \times n$ $= n^K$	$\frac{n^K}{K!}$ not all objects are different
Do not replace	$n \times (n-1) \times$ $\times (n-2) \times \dots \times$ $(n-K+1) =$ $= \frac{n!}{(n-K)!}$	All objects are different \rightarrow $\frac{n!}{(n-K)!} \times \frac{1}{K!} = \binom{n}{K}$

How to solve the problem of K out of n
with replacement but where order
does not matter?

Let's solve $n=2$ problem first:



object 1



object 2

$K=3$

4 possibilities



(1)



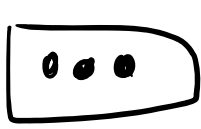
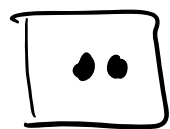
(2)



(3)



(4)



$n=4, K=7$



K dots, $n-1$ box boundaries

$$\binom{K+n-1}{K} = \frac{(K+n-1)!}{K! (n-1)!}$$

ways to distribute

Sampling table

How many ways to choose a sample of k objects out of population of n objects?

	Order matters	Order does not matter
Replacement	$(n)^k$	Difficult: $\binom{n+k-1}{k} = \frac{(n+k-1)!}{(n-1)!k!}$
No replacement	$\frac{n(n-1)(n-2)\dots(n-k+1)}{1} = \frac{n!}{(n-k)!}$	$\binom{n}{k} = \frac{n!}{(n-k)!k!}$

Example

- A DNA of 100 bases is characterized by its numbers of 4 nucleotides:
 d_A , d_C , d_G , and d_T ($d_A + d_C + d_G + d_T = 100$)
- **I don't care about the sequence** (only about the total numbers of A,C,G, and T)
- How many distinct combinations of d_A , d_C , d_G , and d_T are out there?

Probability Axioms,
Conditional Probability,
Statistical (In)dependence,
Circuit Problems

Axioms of probability

Probability is a number that is assigned to each member of a collection of events from a random experiment that satisfies the following properties:

If S is the sample space and E is any event in a random experiment,

(1) $P(S) = 1$

(2) $0 \leq P(E) \leq 1$

(3) For two events E_1 and E_2 with $E_1 \cap E_2 = \emptyset$

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$

$$P(\emptyset) = 0$$

These axioms imply that:

$$P(E') = 1 - P(E)$$

if the event E_1 is contained in the event E_2

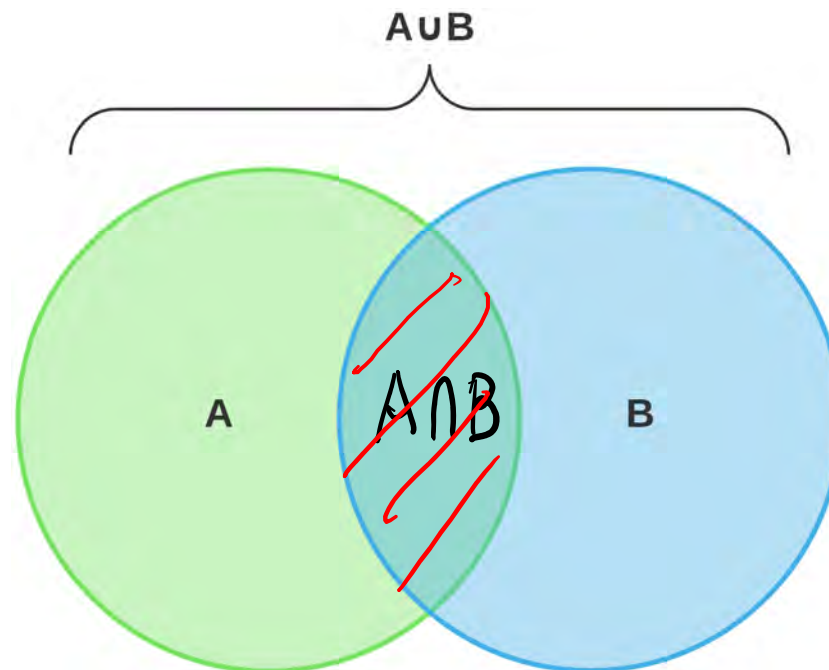
$$P(E_1) \leq P(E_2)$$

Addition rules following from the Axiom (3)

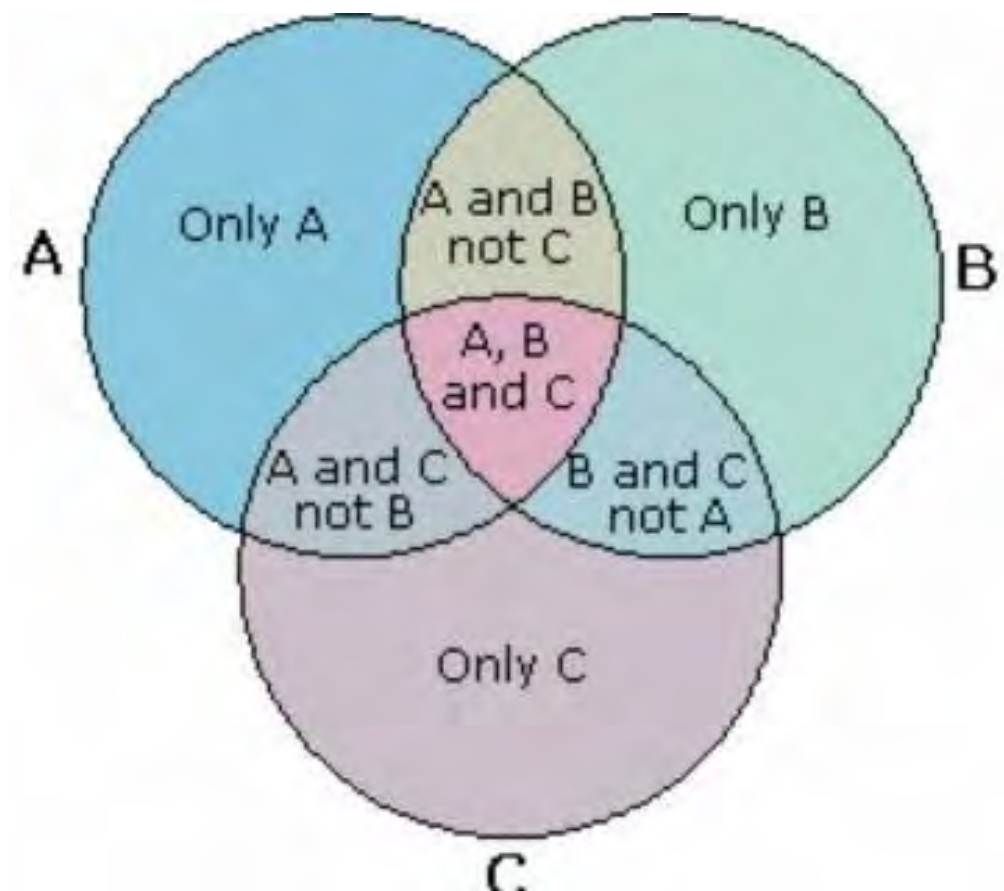
If A and B are mutually exclusive events, i.e. $A \cap B = \emptyset$

$$P(A \cup B) = P(A) + P(B) \quad (2-2)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) \quad (2-1)$$



$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - \\ - P(A \cap B) - P(A \cap C) - P(B \cap C) + \\ + P(A \cap B \cap C).$$



Conditional probability

The **conditional probability** of an event B given an event A , denoted as $P(B|A)$, is

$$P(B|A) = P(A \cap B)/P(A)$$

for $P(A) > 0$.

This definition can be understood in a special case in which all outcomes of a random experiment are equally likely. If there are n total outcomes,

$$P(A) = (\text{number of outcomes in } A)/n$$

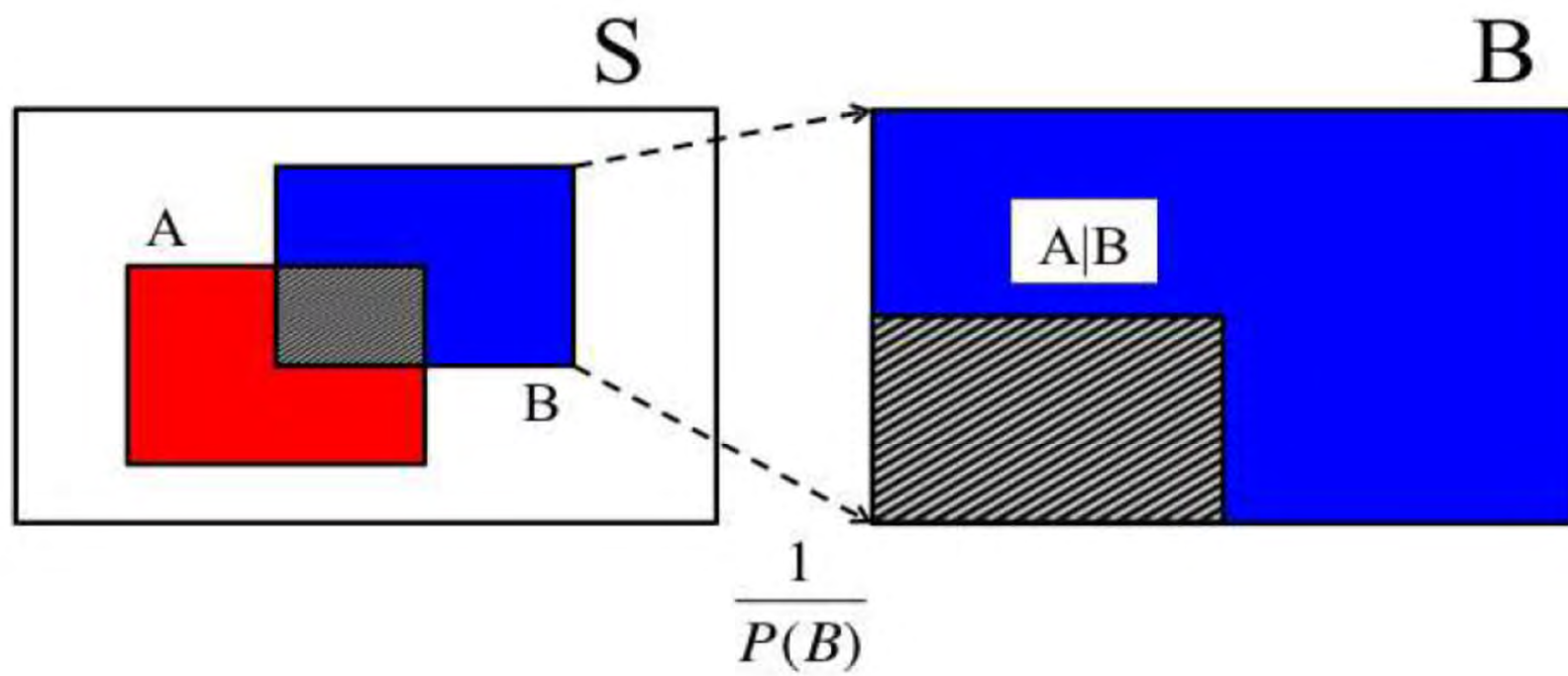
Also,

$$P(A \cap B) = (\text{number of outcomes in } A \cap B)/n$$

Consequently,

$$P(A \cap B)/P(A) = \frac{\text{number of outcomes in } A \cap B}{\text{number of outcomes in } A}$$

Therefore, $P(B|A)$ can be interpreted as the relative frequency of event B among the trials that produce an outcome in event A .



Multiplication rule

is just definition of conditional probability

$$P(\textcolor{red}{B} \mid \textcolor{blue}{A}) = P(\textcolor{red}{B} \cap \textcolor{blue}{A}) / P(\textcolor{blue}{A}) \rightarrow$$

$$P(\textcolor{red}{B} \cap \textcolor{blue}{A}) = P(\textcolor{red}{B} \mid \textcolor{blue}{A}) \cdot P(\textcolor{blue}{A})$$

Drake equation

$$N = R^* \cdot f_p \cdot n_e \cdot f_l \cdot f_i \cdot f_c \cdot L$$

- N = The number of civilizations in The Milky Way Galaxy whose electromagnetic emissions are detectable.
- R^* = The rate of formation of stars suitable for the development of intelligent life.
- f_p = The fraction of those stars with planetary systems.
- n_e = The number of planets, per solar system, with an environment suitable for life.
- f_l = The fraction of suitable planets on which life actually appears.
- f_i = The fraction of life bearing planets on which intelligent life emerges.
- f_c = The fraction of civilizations that develop a technology that releases detectable signs of their existence into space.
- L = The length of time such civilizations release them

Statistically independent events

Always true: $P(A \cap B) = P(A | B) \cdot P(B) = P(B | A) \cdot P(A)$

■ Two events

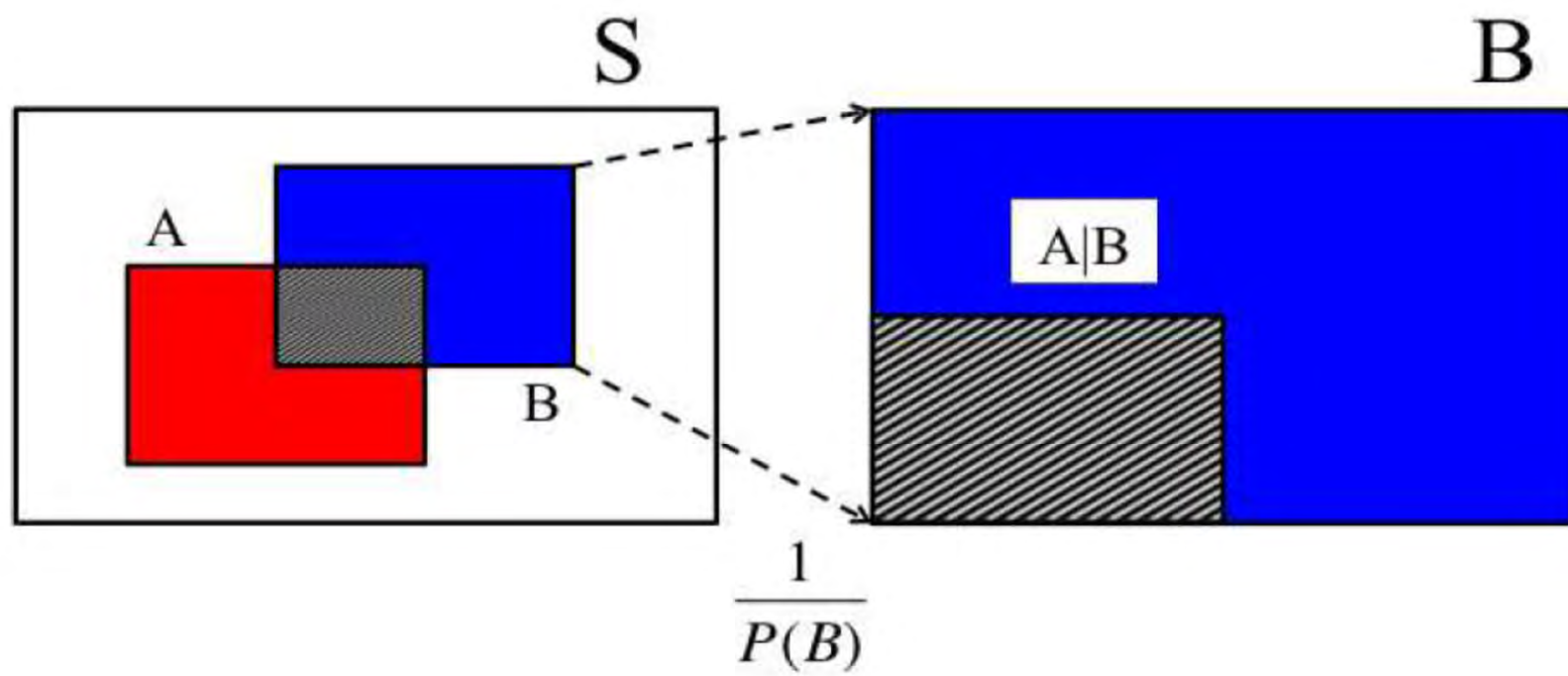
Two events are **independent** if **any one** of the following equivalent statements is true:

- (1) $P(A|B) = P(A)$
- (2) $P(B|A) = P(B)$
- (3) $P(A \cap B) = P(A)P(B)$

■ Multiple events

The events E_1, E_2, \dots, E_n are independent if and only if for any subset of these events $E_{i_1}, E_{i_2}, \dots, E_{i_k}$,

$$P(E_{i_1} \cap E_{i_2} \cap \dots \cap E_{i_k}) = P(E_{i_1}) \times P(E_{i_2}) \times \dots \times P(E_{i_k})$$



Example 3.10. Let an experiment consist of drawing a card at random from a standard deck of 52 playing cards. Define events A and B as “the card is a ♠” and “the card is a queen.” Are the events A and B independent? By definition, $P(A \cdot B) = P(Q\spadesuit) = \frac{1}{52}$. This is the product of $P(\spadesuit) = \frac{13}{52}$ and $P(Q) = \frac{4}{52}$, and events A and B in question are independent. In this situation, intuition provides no help. Now, pretend that the $2\heartsuit$ is drawn and excluded from the deck prior to the experiment. Events A and B become dependent since

$$\mathbb{P}(A) \cdot \mathbb{P}(B) = \frac{13}{51} \cdot \frac{4}{51} \neq \frac{1}{51} = \mathbb{P}(A \cdot B).$$

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED

WHY IS SPACE BLACK

WHY IS OUTER SPACE SO COLD

WHY ARE THERE PYRAMIDS ON THE MOON

WHY IS NASA SHUTTING DOWN

WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS SEX
SO IMPORTANT



WHY ARE THERE
GHOSTS



WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS

WHY AREN'T MY QUAIL EGGS HATCHING

WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY IS THERE AN OWL IN MY BACKYARD

WHY IS THERE AN OWL OUTSIDE MY WINDOW

WHY IS THERE AN OWL ON THE DOLLAR BILL

WHY DO OWLS ATTACK PEOPLE

WHY ARE AK 47s SO EXPENSIVE

WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE

WHY ARE THERE GODS

WHY ARE THERE TWO SPOCKS

WHY ARE CIGARETTES LEGAL

WHY ARE THERE DUCKS IN MY POOL

WHY IS JESUS WHITE

WHY IS THERE LIQUID IN MY EAR

WHY DO Q TIPS FEEL GOOD

WHY DO GOOD PEOPLE DIE

WHY AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT

WHY ARE ULTRASOUND MACHINES EXPENSIVE

WHY IS STEALING WRONG

WHY ARE THERE FIREWORKS

WHY ARE THERE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS

WHY DO I SAY UH

WHY IS SEA SALT BETTER

WHY ARE THERE TREES IN THE MIDDLE OF FIELDS

WHY IS THERE NOT A POKEMON MMO

WHY IS THERE LAUGHING IN TV SHOWS

WHY ARE THERE DOORS ON THE FREEWAY

WHY ARE THERE SO MANY SVCHOST.EXE RUNNING

WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA

WHY ARE THERE SCARY SOUNDS IN MINECRAFT

WHY IS THERE KICKING IN MY STOMACH

WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES

WHY DO SNAKES EXIST

WHY DO OYSTERS HAVE PEARLS

WHY ARE DUCKS CALLED DUCKS

WHY DO THEY CALL IT THE CLAP

WHY ARE KYLE AND CARTMAN FRIENDS

WHY IS THERE AN ARROW ON AANG'S HEAD

WHY ARE TEXT MESSAGES BLUE

WHY ARE THERE MUSTACHES ON CLOTHES

WHY ARE THERE MUSTACHES ON CARS

WHY ARE THERE MUSTACHES EVERYWHERE

WHY ARE THERE SO MANY BIRDS IN OHIO

WHY IS THERE SO MUCH RAIN IN OHIO

WHY IS OHIO WEATHER SO WEIRD

WHY ARE THERE BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP

WHY AREN'T THERE VARICOSE ARTERIES

WHY ARE OLD KINGDOMS DIFFERENT

WHY ARE THERE SQUIRRELS



WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL

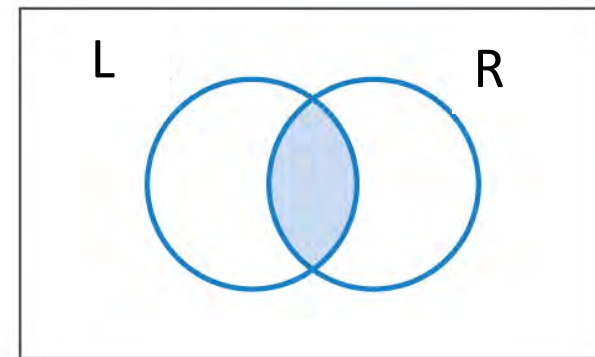
WHY IS THERE HELL IF GOD FORGIVES

WHY IS THERE NO GPS IN LAPTOPS
WHY DO KNEES CLICK
WHY AREN'T THERE E GRADES
WHY IS ISOLATION BAD
WHY DO BOYS LIKE ME
WHY DON'T BOYS LIKE ME
WHY IS THERE ALWAYS A JAVA UPDATE
WHY ARE THERE RED DOTS ON MY THIGHS
WHY IS LYING GOOD

WHY IS GPS FREE

Series Circuit

This circuit operates only if there is **at least one path of functional devices** from left to right. The **probability** that **each device functions** is shown on the graph. Assume that the **devices fail independently**. What is the probability that the circuit operates?

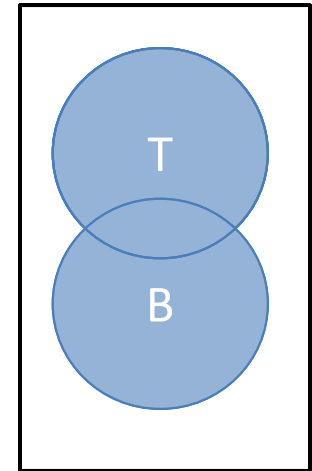
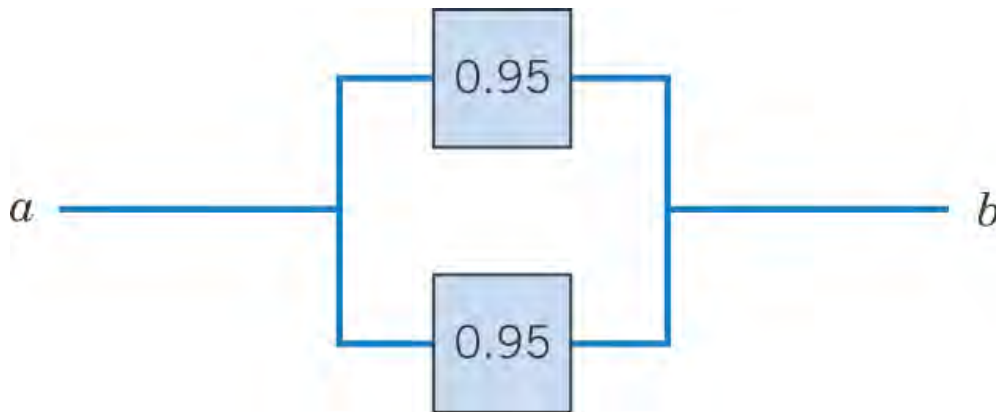


Let L & R denote the events that the left and right devices operate. The probability that the circuit operates is:

$$P(L \text{ and } R) = P(L \cap R) = P(L) * P(R) = 0.8 * 0.9 = 0.72.$$

Parallel Circuit

This circuit operates only if there is a path of functional devices from left to right. The probability that each device functions is shown. Each device fails independently.

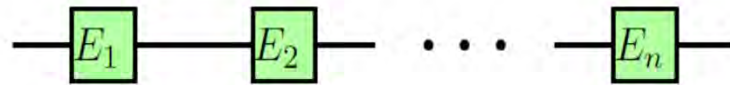


Let T & B denote the events that the top and bottom devices operate. The probability that the circuit operates is:

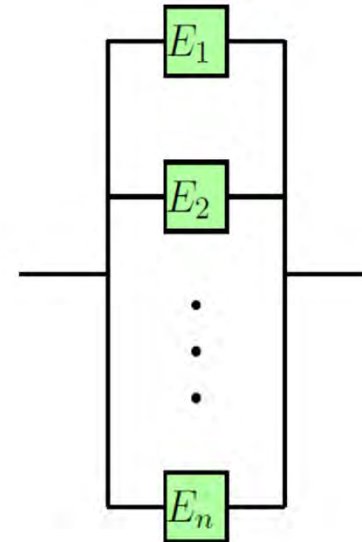
$$P(T \cup B) = 1 - P(T' \cap B') = 1 - P(T') * P(B') = 1 - 0.05^2 = 1 - 0.0025 = 0.9975.$$

Duality between parallel and series circuits

$$q_i = 1 - p_i.$$



(a)

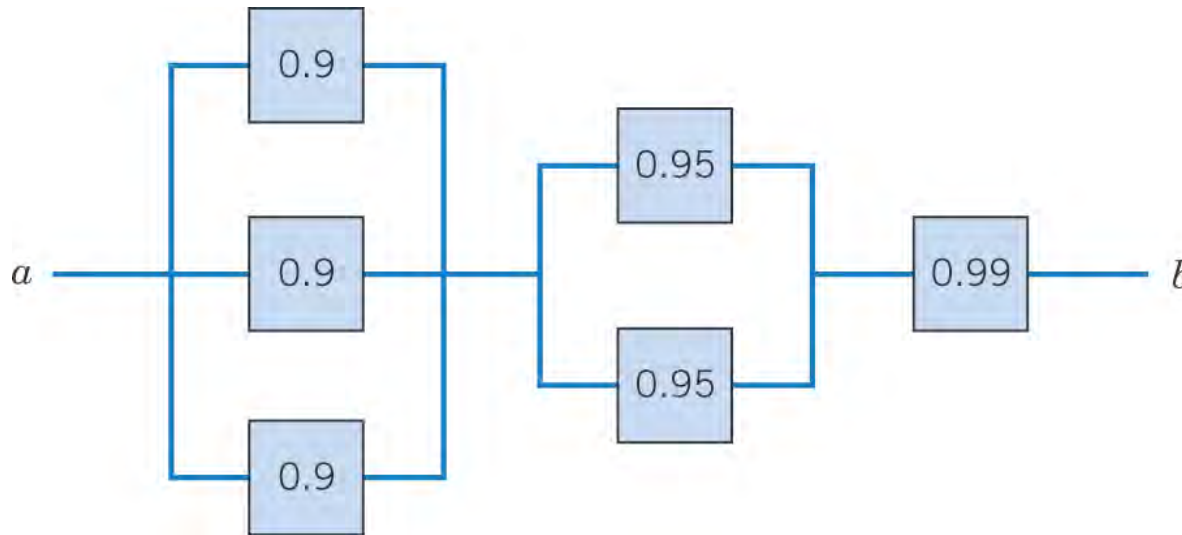


(b)

Connection	Notation	Works with prob	Fails with prob
Serial	$E_1 \cap E_2 \cap \dots \cap E_n$	$p_1 p_2 \dots p_n$	$1 - p_1 p_2 \dots p_n$
Parallel	$E_1 \cup E_2 \cup \dots \cup E_n$	$1 - q_1 q_2 \dots q_n$	$q_1 q_2 \dots q_n$

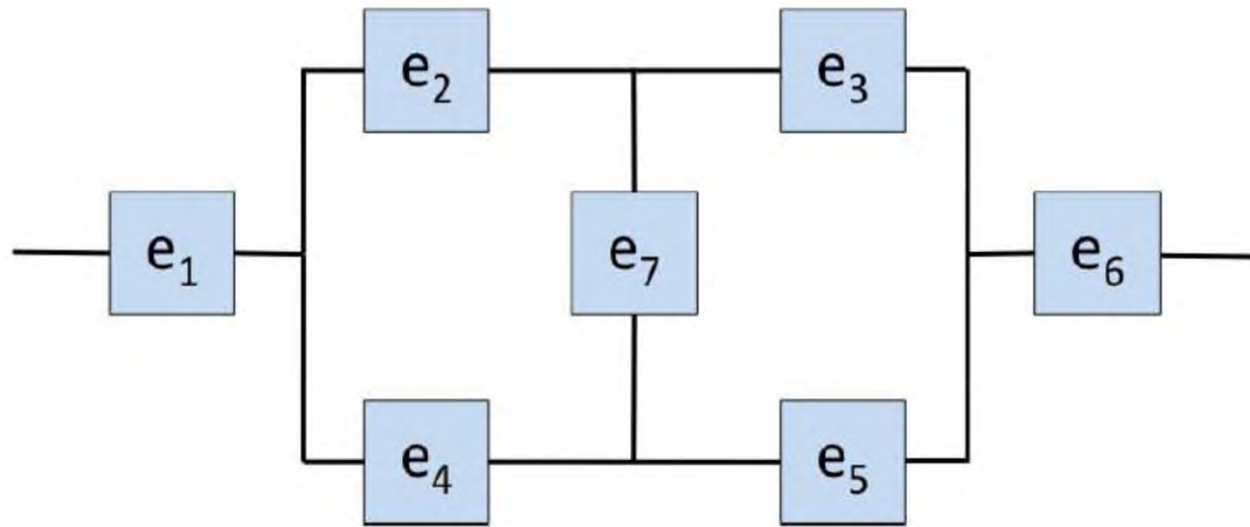
Advanced Circuit

This circuit operates only if there is a path of functional devices from left to right. The probability that each device functions is shown. Each device fails independently.

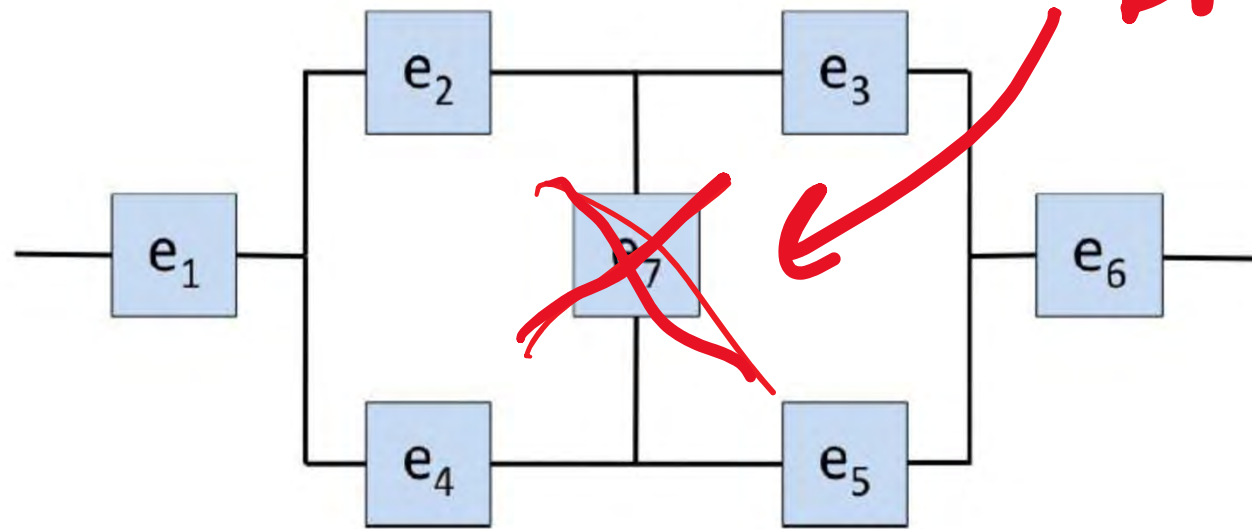


Partition the graph into 3 columns with L & M denoting the left & middle columns.

$P(L) = 1 - 0.1^3$, and $P(M) = 1 - 0.05^2$, so the probability that the circuit operates is: $(1 - 0.1^3)(1 - 0.05^2)(0.99) = 0.9875$ (this is a series of parallel circuits).



Component	e_1	e_2	e_3	e_4	e_5	e_6	e_7
Probability of component working	0.3	0.8	0.2	0.2	0.5	0.6	0.4



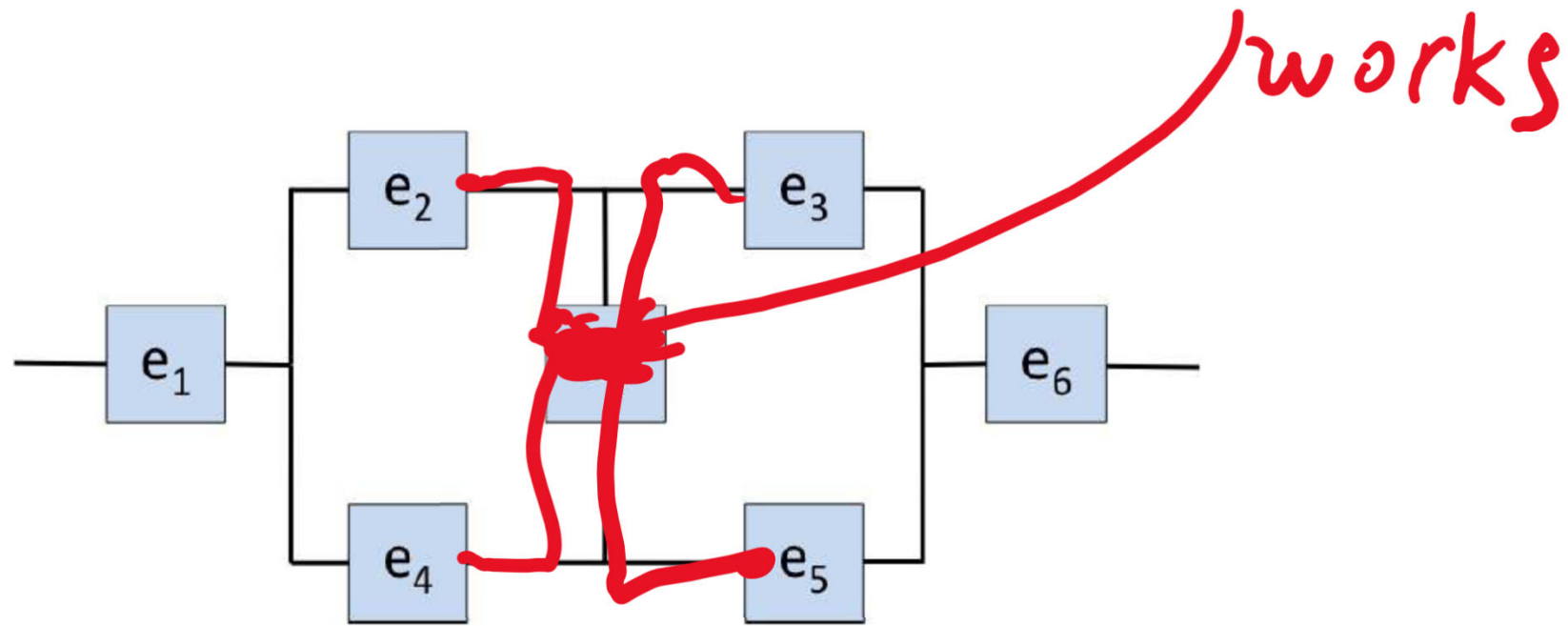
broken

Component	e_1	e_2	e_3	e_4	e_5	e_6	e_7
Probability of component working	0.3	0.8	0.2	0.2	0.5	0.6	0.4

$$\begin{aligned}
 &P(\text{circuit works} \mid e_7 \text{ is broken}) = P(e_1 \text{ works}) * \\
 &[1 - (1 - P(e_2 \text{ works}) * P(e_3 \text{ works})) * (1 - P(e_4 \text{ works}) * P(e_5 \text{ works}))] * \\
 &P(e_6 \text{ works}) = 0.3 * (1 - (1 - 0.8 * 0.2) * (1 - 0.2 * 0.5)) * 0.6 = 0.0439
 \end{aligned}$$

The contribution to total probability:

$$P(\text{circuit works} \mid e_7 \text{ is broken}) * P(e_7 \text{ is broken}) = 0.6 * 0.0439 = 0.0264$$



Component	e_1	e_2	e_3	e_4	e_5	e_6	e_7
Probability of component working	0.3	0.8	0.2	0.2	0.5	0.6	0.4

$P(\text{circuit works} \mid e_7 \text{ works}) = P(e_1 \text{ works}) \cdot$

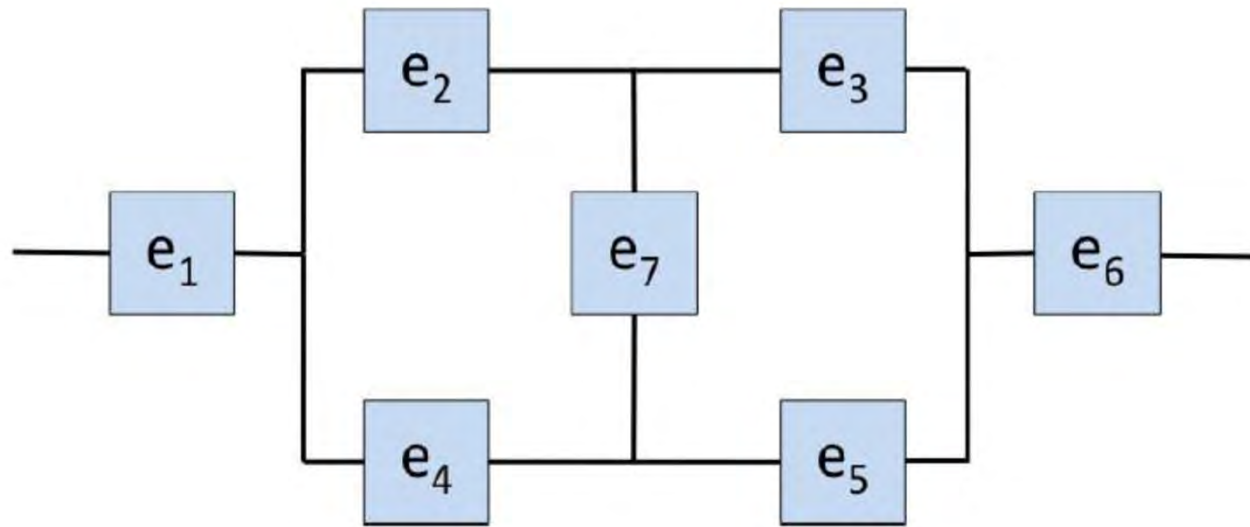
$[1 - (1 - P(e_2 \text{ works})) \cdot (1 - P(e_4 \text{ works}))]$

$\cdot [1 - (1 - P(e_3 \text{ works})) \cdot (1 - P(e_5 \text{ works}))] \cdot$

$P(e_6 \text{ works}) = 0.3 \cdot (1 - (1 - 0.8) \cdot (1 - 0.2)) \cdot (1 - (1 - 0.2) \cdot (1 - 0.5)) \cdot 0.6 = 0.0907$

The contribution to total probability:

$P(\text{circuit works} \mid e_7 \text{ works}) \cdot P(e_7 \text{ works}) = 0.4 \cdot 0.0907 = 0.0363$

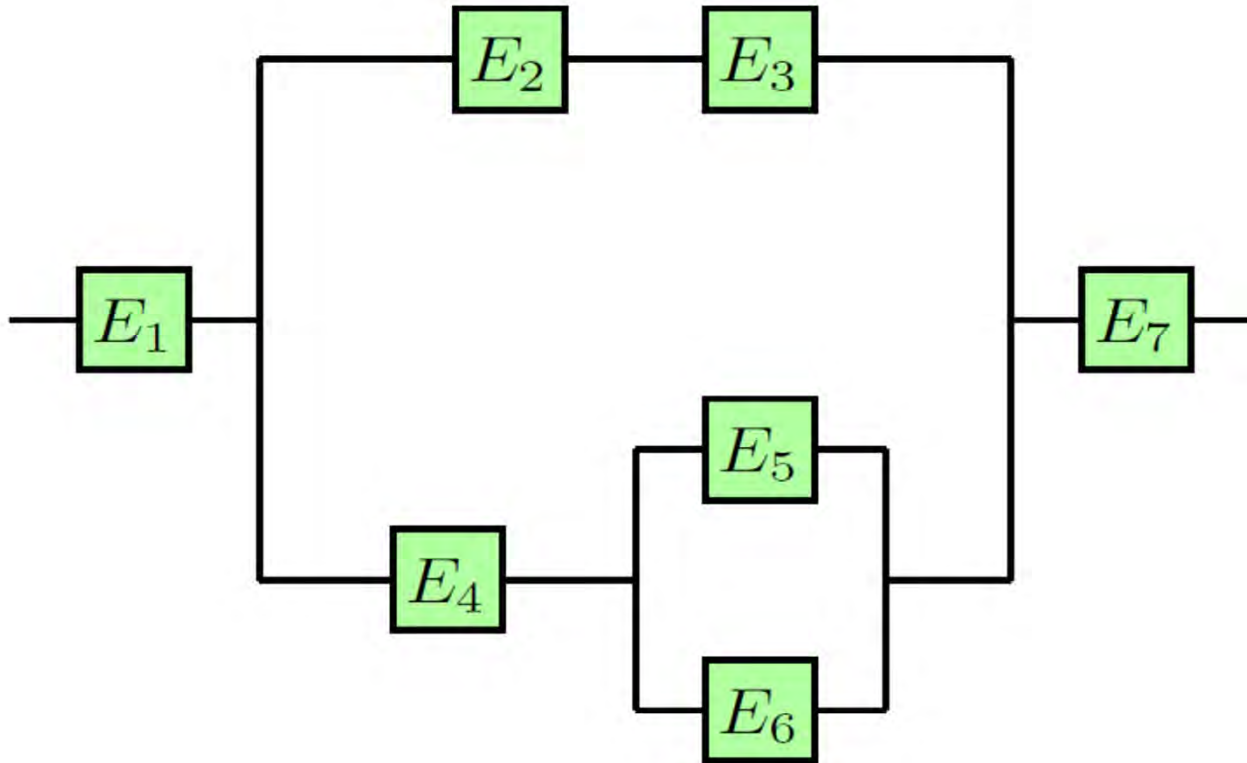


Component	e_1	e_2	e_3	e_4	e_5	e_6	e_7
Probability of component working	0.3	0.8	0.2	0.2	0.5	0.6	0.4

$$\begin{aligned}
 &P(\text{circuit works}) = \\
 &P(\text{circuit works} \mid e_7 \text{ works}) * P(e_7 \text{ works}) + \\
 &P(\text{circuit works} \mid e_7 \text{ is broken}) * P(e_7 \text{ is broken}) = \\
 &= 0.0264 + 0.0363 = 0.0627
 \end{aligned}$$

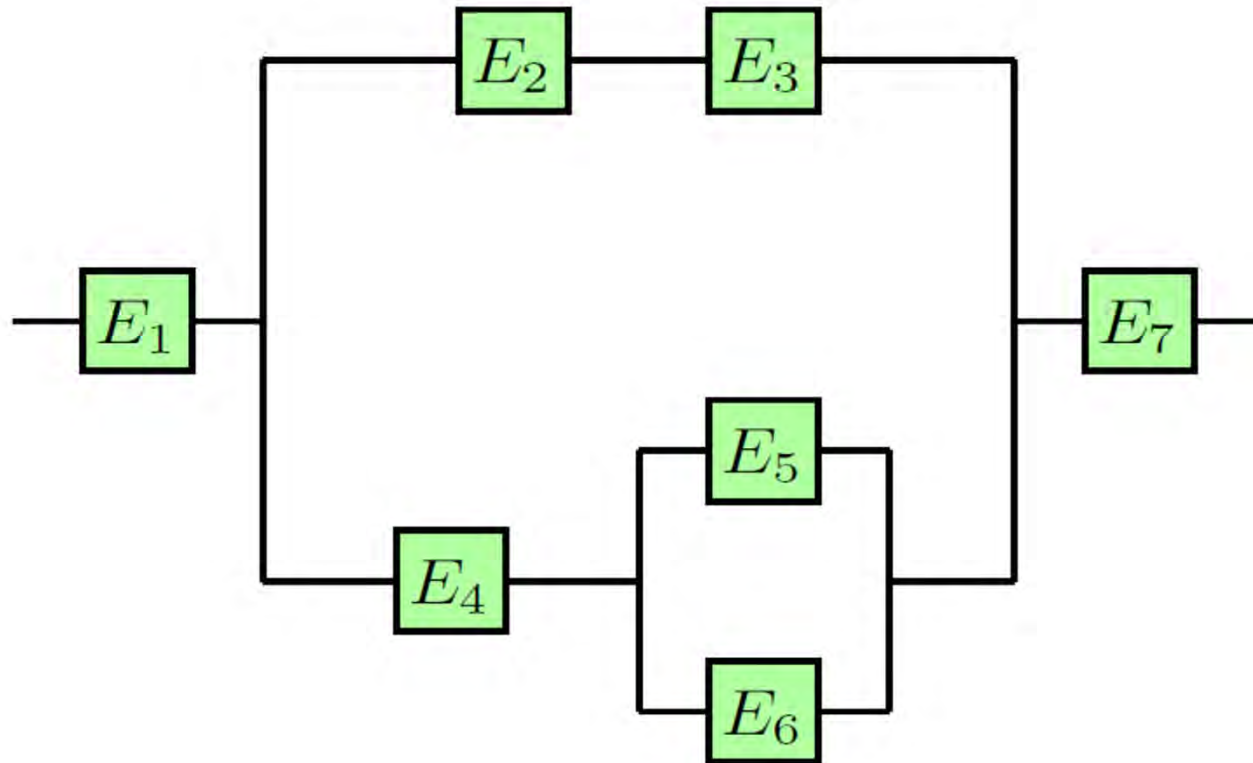
Answer: 6.27%

Circuit \rightarrow Set equation



Component	E_1	E_2	E_3	E_4	E_5	E_6	E_7
Probability of functioning well	0.9	0.5	0.3	0.1	0.4	0.5	0.8

Circuit → Set equation



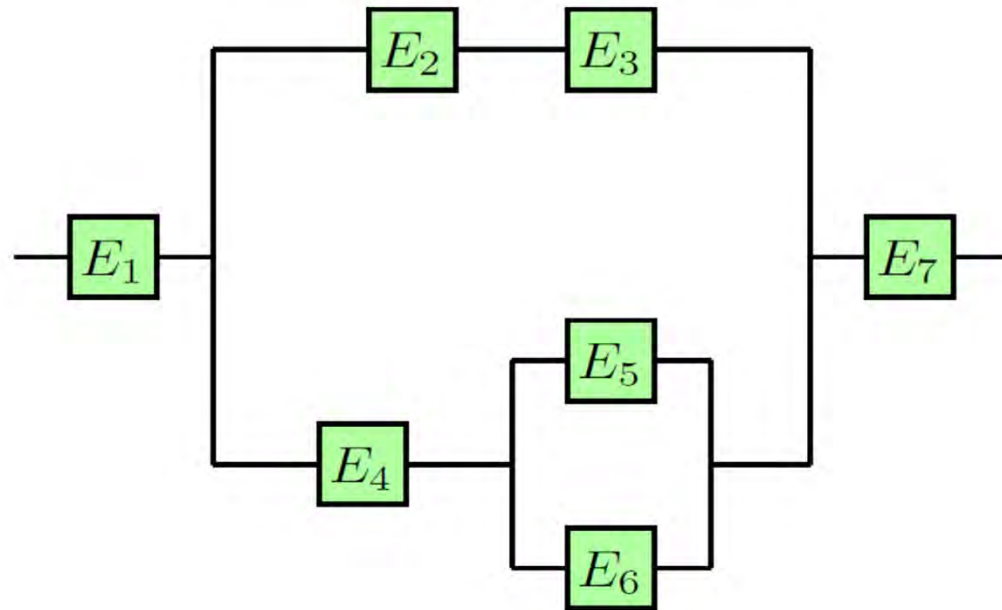
Component	E_1	E_2	E_3	E_4	E_5	E_6	E_7
Probability of functioning well	0.9	0.5	0.3	0.1	0.4	0.5	0.8

$$E_1 \cap [(E_2 \cap E_3) \cup (E_4 \cap (E_5 \cup E_6))] \cap E_7.$$

$$P(\text{Works}) = 0.9 \cdot (1 - (1 - 0.5 \cdot 0.3) \cdot (1 - 0.1 \cdot (1 - 0.6 \cdot 0.5))) \cdot 0.8 = 0.15084$$

Matlab group exercise

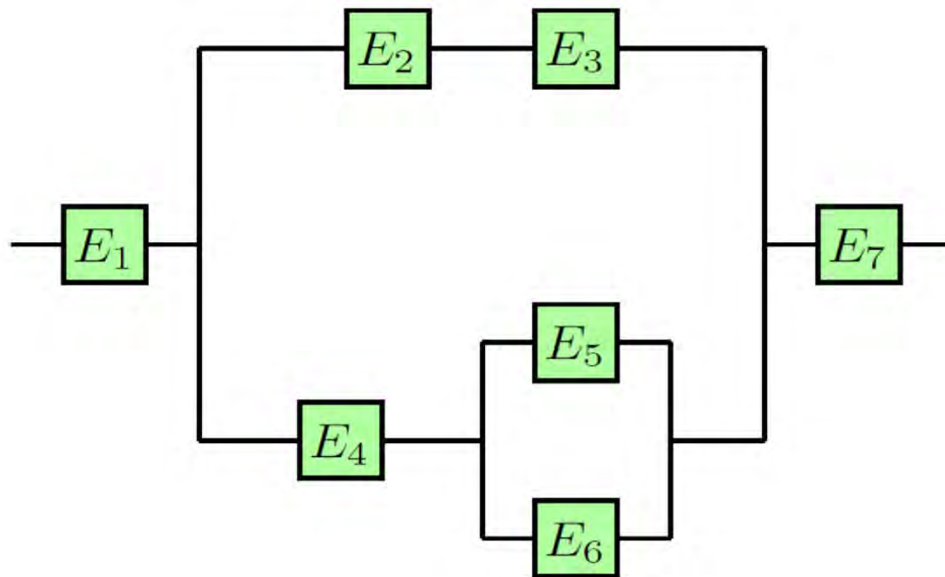
- Test our result for this circuit.
- Use circuit_template.m on the website



Component	E_1	E_2	E_3	E_4	E_5	E_6	E_7
Probability of functioning well	0.9	0.5	0.3	0.1	0.4	0.5	0.8

Matlab group exercise

- Test our result for this circuit.
- Download `circuit_template.m` from the website



Component	E_1	E_2	E_3	E_4	E_5	E_6	E_7
Probability of functioning well	0.9	0.5	0.3	0.1	0.4	0.5	0.8

$$P(\text{Works}) = 0.9 \cdot (1 - (1 - 0.5 \cdot 0.3) \cdot (1 - 0.1 \cdot (1 - 0.6 \cdot 0.5))) \cdot 0.8 = 0.15084$$

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED

WHY IS SPACE BLACK

WHY IS OUTER SPACE SO COLD

WHY ARE THERE PYRAMIDS ON THE MOON

WHY IS NASA SHUTTING DOWN

WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS ISOLATION BAD

WHY DO BOYS LIKE ME

WHY DON'T BOYS LIKE ME

WHY IS THERE ALWAYS A JAVA UPDATE

WHY ARE THERE RED DOTS ON MY THIGHS

WHY IS LYING GOOD

WHY ARE THERE
GHOSTS



WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS

WHY AREN'T MY QUAIL EGGS HATCHING

WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE

WHY AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG

WHY ARE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS

WHY DO I SAY UH

WHY IS SEA SALT BETTER

WHY ARE THERE TREES IN THE MIDDLE OF FIELDS

WHY IS THERE NOT A POKEMON MMO

WHY IS THERE LAUGHING IN TV SHOWS

WHY ARE THERE DOORS ON THE FREEWAY

WHY ARE THERE SO MANY SVCHOST.EXE RUNNING

WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA

WHY ARE THERE SCARY SOUNDS IN MINECRAFT

WHY IS THERE KICKING IN MY STOMACH

WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES

WHY DO SNAKES EXIST

WHY DO OYSTERS HAVE PEARLS

WHY ARE DUCKS CALLED DUCKS

WHY DO THEY CALL IT THE CLAP

WHY ARE KYLE AND CARTMAN FRIENDS

WHY IS THERE AN ARROW ON AANG'S HEAD

WHY ARE TEXT MESSAGES BLUE

WHY ARE THERE MUSTACHES ON CLOTHES

WHY ARE THERE MUSTACHES ON CARS

WHY ARE THERE MUSTACHES EVERYWHERE

WHY ARE THERE SO MANY BIRDS IN OHIO

WHY IS THERE SO MUCH RAIN IN OHIO

WHY IS OHIO WEATHER SO WEIRD

WHY ARE THERE BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP

WHY AREN'T THERE VARICOSE ARTERIES

WHY ARE OLD KINGDOMS DIFFERENT

WHY ARE THERE SQUIRRELS

WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL



WHY IS SEX
SO IMPORTANT

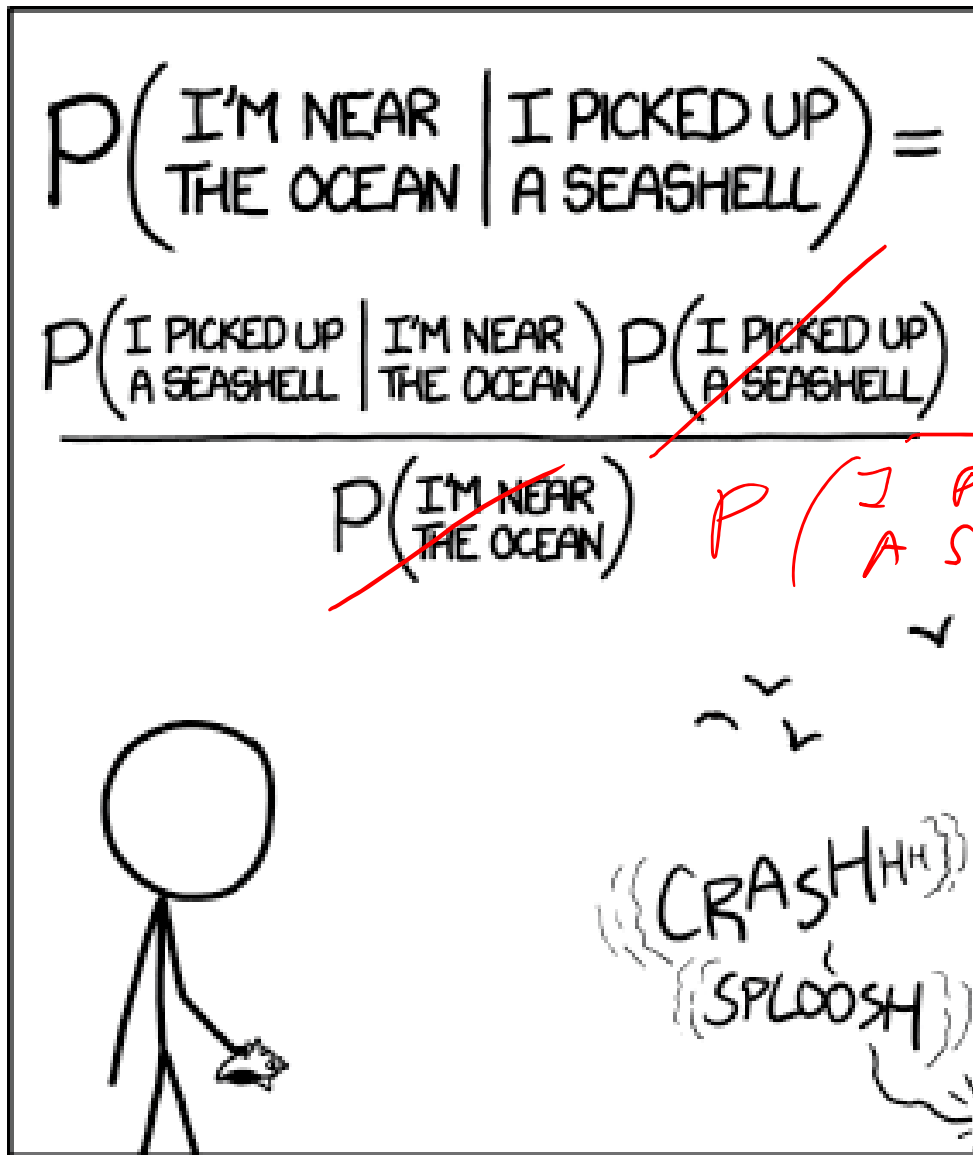


WHY AREN'T THERE DINOSAUR GHOSTS

WHY ARE THERE FEMALE MR NIMES

WHY IS GPS FREE

Reminder:
Conditional probability



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

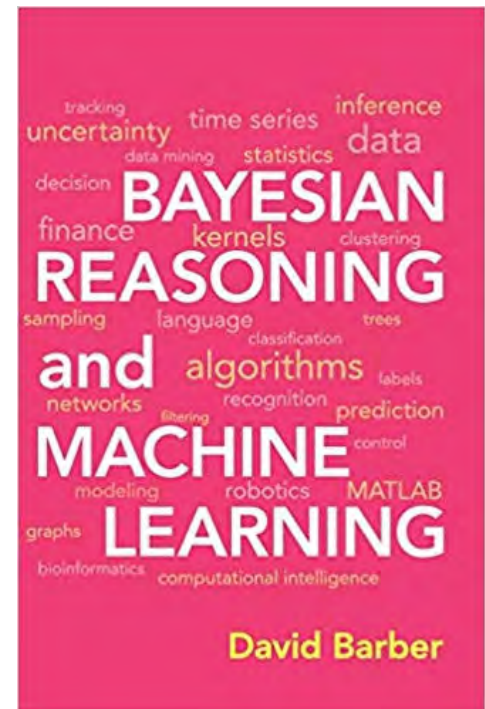
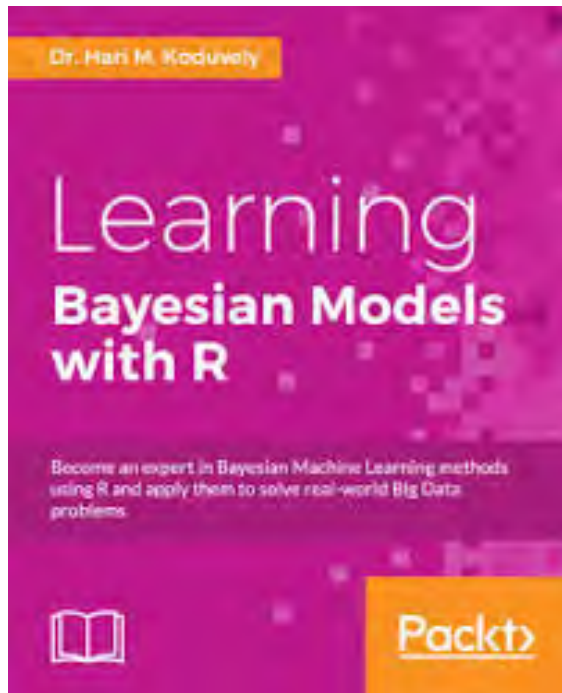
What is wrong in this comics?

If you are not yet reading XKCD comics

<https://xkcd.com/>
you should start

Bayes Theorem

Bayes' theorem



Thomas Bayes (1701-1761)

English statistician, philosopher, and Presbyterian minister

Bayes' theorem was presented in "An Essay towards solving a Problem in the Doctrine of Chances" which was read to the Royal Society in 1763 already after Bayes' death.

Bayes' theorem (simple)

$$P(A \cap B) = \underline{P(A|B)P(B)} = P(B \cap A) = \underline{P(B|A)P(A)}$$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- In Science **we often want to know:**
“**How much faith** should I put into **hypothesis, given the data?**”
or $P(H|D)$ (see also the inductive definition of probability)
- What **we usually can calculate** if the hypothesis/model is OK:
“Assuming that this **hypothesis is true**, what is the **probability of the observed data?**” or $P(D|H)$
- Bayes' theorem can help: $P(H|D) = P(D|H) \cdot P(H) / P(D)$
- The problem is $P(H)$ (so-called prior) is often **not known**

Bayes' theorem (continued)

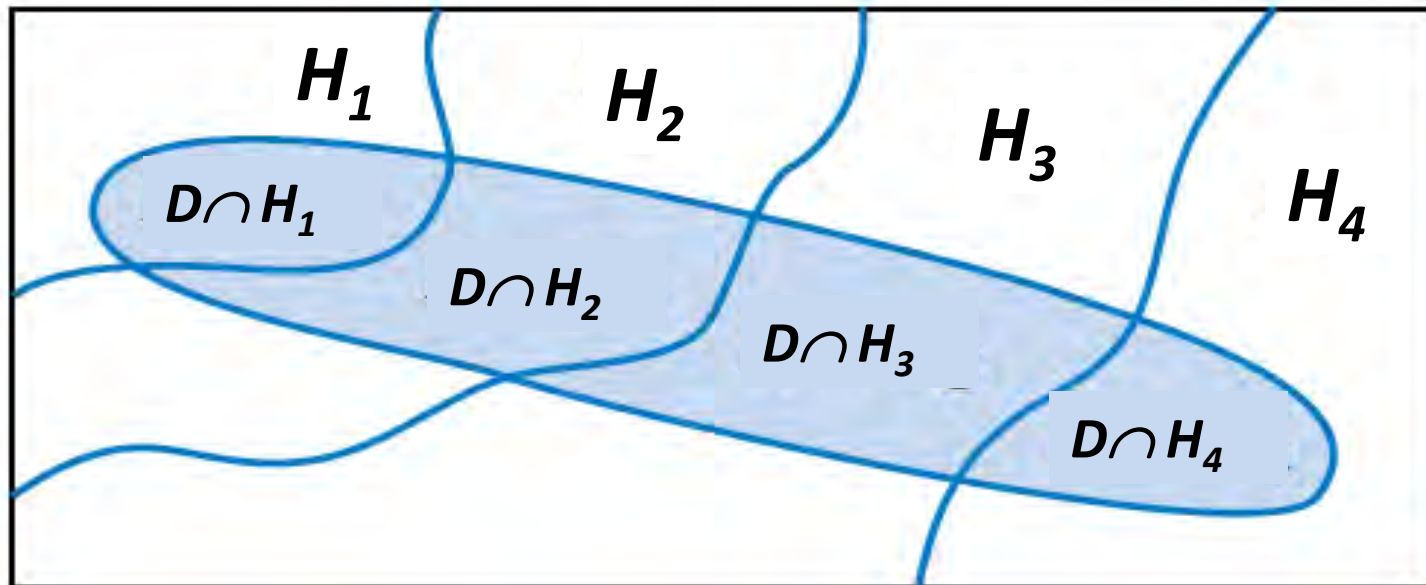
Works best with **exhaustive** and **mutually-exclusive** hypotheses:

H_1, H_2, \dots, H_n such that $H_1 \cup H_2 \cup H_3 \dots \cup H_n = S$ and $H_i \cap H_j = \emptyset$ for $i \neq j$

$$P(H_k|D) = P(D|H_k) \cdot P(H_k) / P(D)$$

where:

$$P(D) = P(D|H_1) \cdot P(H_1) + P(D|H_2) \cdot P(H_2) + \dots + P(D|H_n) \cdot P(H_n)$$



Secretary problem

- An **employer** has a known number – n – of **applicants** for a secretary position, whom are **interviewed one at a time**
- Employer can easily **evaluate and rank** applicants relative to each other but has no idea of the overall distribution of their quality
- Employer has only one chance to choose the secretary: gives **yes/no answer in the end of each interview** and cannot go back to rejected applicants
- How can employer **maximize the probability to choose the best secretary** among all applicants?



Martin Gardner (1914 – 2010)
Described the “secretary problem”
in *Scientific American* 1960.

was an American popular
mathematics and popular
science writer. Best known
for “recreational mathematics”:
He was behind the
“Mathematical Games” section
in *Scientific American*.



Eugene Dynkin (1924 – 2014)
solved this problem in 1963.
He referred to it as a “picky bride
problem”

was a Soviet and later American
mathematician, member of the
US National Academy of Science.
He has made contributions to the
fields of probability and algebra.
The Dynkin diagram, the Dynkin
system, and Dynkin's lemma are
all named after him.

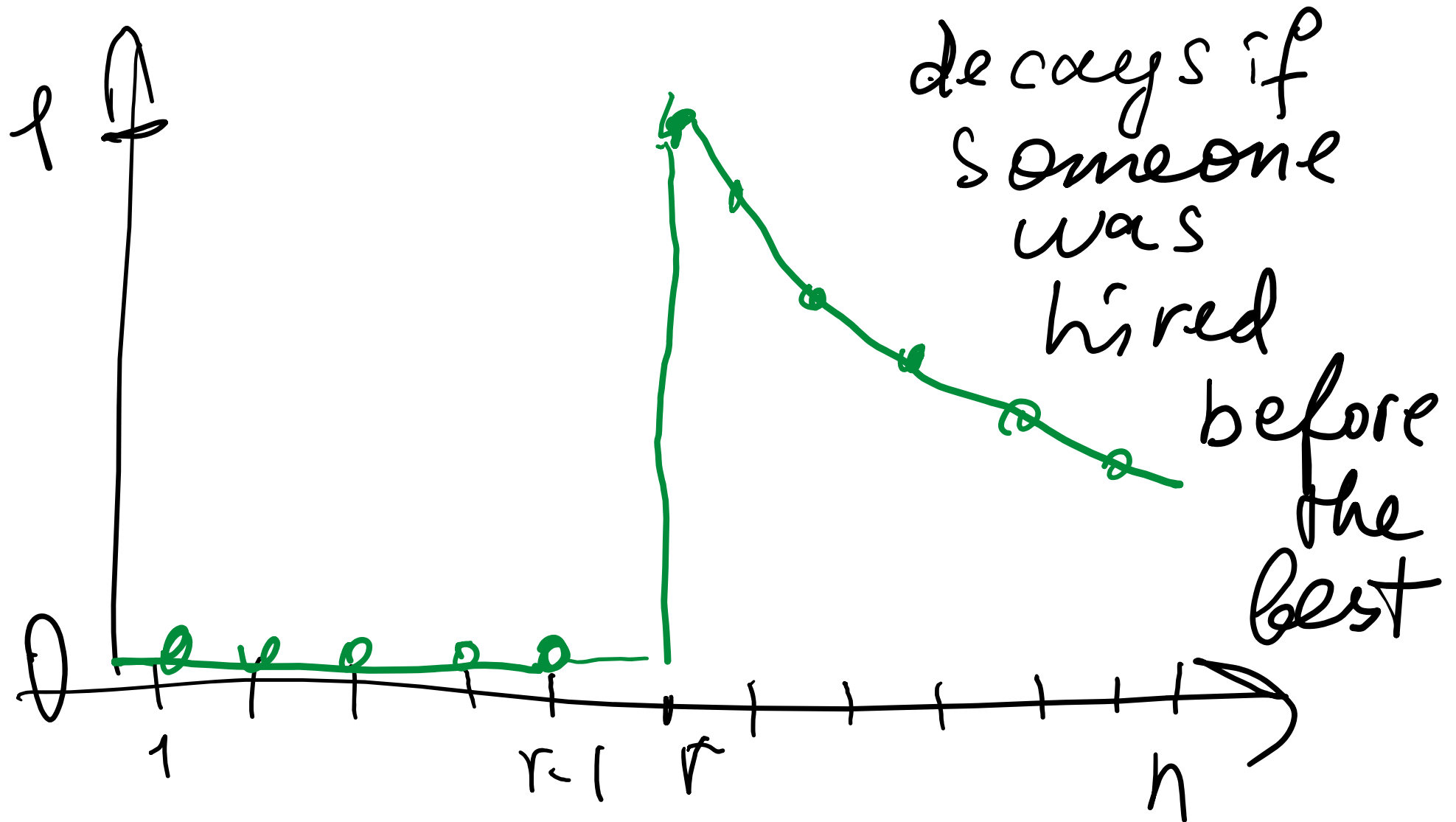
Who solved the secretary problem?

- Gardner outlined the solution in Sci Am 1960 but gave no formal proof
- Solution by Lindey was published in 1961:
Lindey, D. V. (1961). Dynamic programming and decision theory. Appl. Statist. 10 39-51
- Dynkin's paper was published in 1963:
Dynkin, E. B. (1963). The optimum choice of the instant for stopping a Markov process. Soviet Math. Dokl. 4 627-629
- When the celebrated German astronomer, Johannes Kepler (1571-1630), lost his first wife to cholera in 1611, he set about finding a new wife
- He spent 2 years on the process, had 11 candidates and married the 5th candidate ($11/e \sim 4$ so he married the first after)

What should the employer do?

- Employer does not know the distribution of the quality of applicants and has to learn it on the fly
- Algorithm: look at the first $r-1$ applicants, remember the best among them
- Hire the first among next $n-r+1$ applicants who is better than the best among the first r applicants
- How to choose r ?
- When r is too small – not enough information: the best among r is not very good. You are likely to hire a bad secretary
- When r is too large (e.g. $r=n-1$) – you procrastinated for too long! You have almost all the information, but you will have to hire the last applicant who is (likely) not particularly good

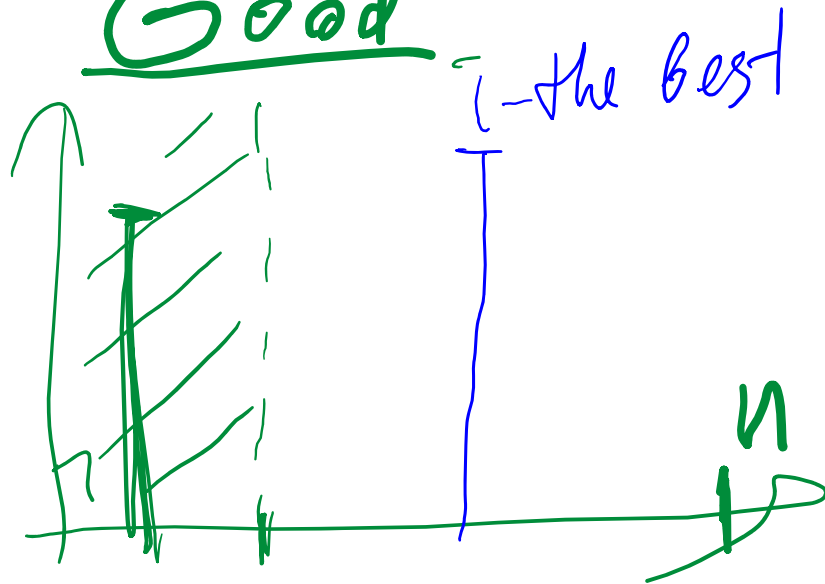
Probability of hiring the best candidate
if he/she has # i in the queue



Look at $i-1$ candidates
before the best

$$\text{Prob} = \frac{r-1}{i-1}$$

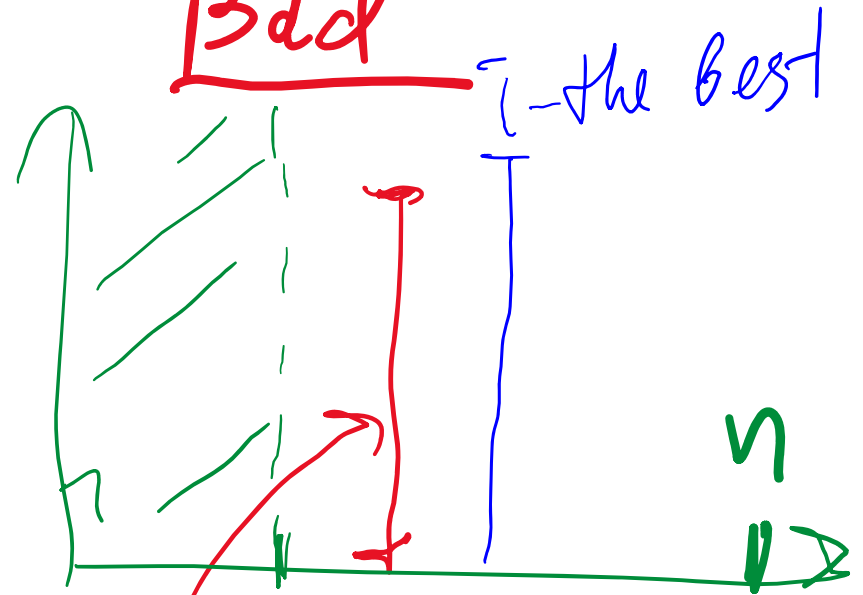
Good



the best among $i-1$

$$\text{Prob} = \frac{i-r}{i-1}$$

Bad



the best among $i-1$

$$\begin{aligned}
P(r) &= \sum_{i=1}^n P(\text{applicant } i \text{ is selected} \cap \text{applicant } i \text{ is the best}) \\
&= \sum_{i=1}^n P(\text{applicant } i \text{ is selected} | \text{applicant } i \text{ is the best}) \times P(\text{applicant } i \text{ is the best}) \\
&= \left[\sum_{i=1}^{r-1} 0 + \sum_{i=r}^n P \left(\begin{array}{c} \text{the best of the first } i-1 \text{ applicants} \\ \text{is in the first } r-1 \text{ applicants} \end{array} \middle| \text{applicant } i \text{ is the best} \right) \right] \times \frac{1}{n} \\
&= \sum_{i=r}^n \frac{r-1}{i-1} \times \frac{1}{n} = \frac{r-1}{n} \sum_{i=r}^n \frac{1}{i-1}.
\end{aligned}$$

$$P(r) = \frac{r-1}{n} \sum_{i=r}^n \frac{1}{i-1}.$$

Letting n tend to infinity, writing x as the limit of r/n , using t for i/n and dt for $1/n$,

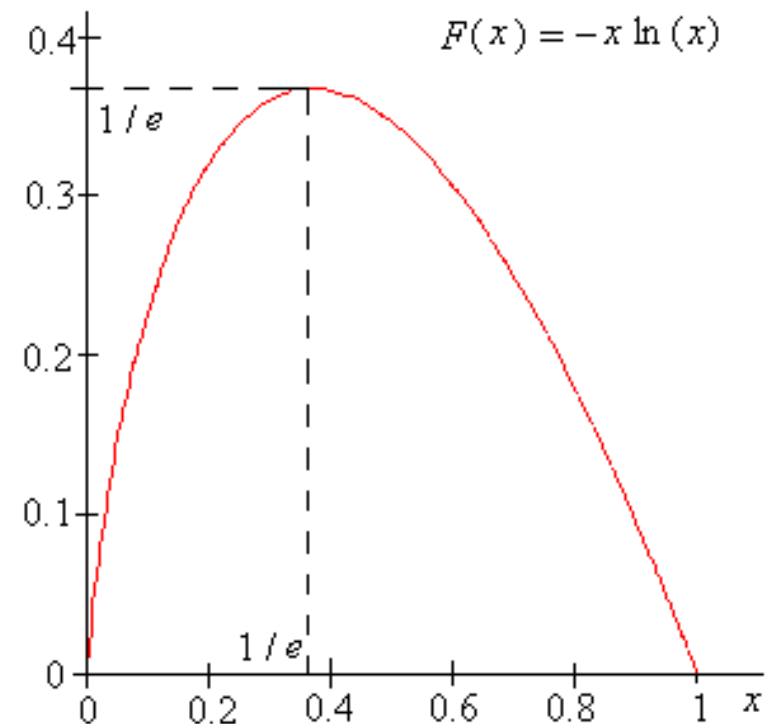
$$P(x) = x \int_x^1 \frac{1}{t} dt = -x \ln(x).$$

$$dP(x)/dx = -\ln(x) - 1$$

$$-\ln(x^*) - 1 = 0$$

$$x^* = 1/e = 0.3679$$

Probability of picking the best applicant is also $1/e = 0.3679$



Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED

WHY IS SPACE BLACK

WHY IS OUTER SPACE SO COLD

WHY ARE THERE PYRAMIDS ON THE MOON

WHY IS NASA SHUTTING DOWN

WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS SEX
SO IMPORTANT



WHY ARE THERE
GHOSTS



WHY IS THERE AN OWL IN MY BACKYARD

WHY IS THERE AN OWL OUTSIDE MY WINDOW

WHY IS THERE AN OWL ON THE DOLLAR BILL

WHY DO OWLS ATTACK PEOPLE

WHY ARE AK 47s SO EXPENSIVE

WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE

WHY ARE THERE GODS

WHY ARE THERE TWO SPOCKS

WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS

WHY AREN'T MY QUAIL EGGS HATCHING

WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE

WHY AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG

WHY ARE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS
WHY DO I SAY UH
WHY IS SEA SALT BETTER

WHY ARE THERE TREES IN THE MIDDLE OF FIELDS
WHY IS THERE NOT A POKEMON MMO
WHY IS THERE LAUGHING IN TV SHOWS
WHY ARE THERE DOORS ON THE FREEWAY
WHY ARE THERE SO MANY SVCHOST.EXE RUNNING
WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA
WHY ARE THERE SCARY SOUNDS IN MINECRAFT
WHY IS THERE KICKING IN MY STOMACH
WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES
WHY DO SNAKES EXIST
WHY DO OYSTERS HAVE PEARLS
WHY ARE DUCKS CALLED DUCKS
WHY DO THEY CALL IT THE CLAP
WHY ARE KYLE AND CARTMAN FRIENDS
WHY IS THERE AN ARROW ON AANG'S HEAD
WHY ARE TEXT MESSAGES BLUE
WHY ARE THERE MUSTACHES ON CLOTHES
WHY ARE THERE MUSTACHES ON CARS
WHY ARE THERE MUSTACHES EVERYWHERE
WHY ARE THERE SO MANY BIRDS IN OHIO
WHY IS THERE SO MUCH RAIN IN OHIO
WHY IS OHIO WEATHER SO WEIRD

WHY ARE THERE BRIDESMAIDS
WHY DO DYING PEOPLE REACH UP
WHY AREN'T THERE VARICOSE ARTERIES
WHY ARE OLD KINGDOMS DIFFERENT

WHY ARE THERE
SQUIRRELS



WHY IS PROGRAMMING SO HARD
WHY IS THERE A 0 OHM RESISTOR
WHY DO AMERICANS HATE SOCCER
WHY DO RHYMES SOUND GOOD
WHY DO TREES DIE
WHY IS THERE NO SOUND ON CNN
WHY AREN'T POKEMON REAL
WHY AREN'T BULLETS SHARP
WHY DO DREAMS SEEM SO REAL

WHY IS THERE HELL IF GOD FORGIVES
WHY IS GPS FREE

WHY ARE THERE TINY SPIDERS IN MY HOUSE
WHY DO SPIDERS COME INSIDE
WHY ARE THERE HUGE SPIDERS IN MY HOUSE
WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE
WHY ARE THERE SPIDERS IN MY ROOM
WHY ARE THERE SO MANY SPIDERS IN MY ROOM
WHY DO SPIDER BITES ITCH
WHY IS DYING SO SCARY
WHY IS THERE NO GPS IN LAPTOPS
WHY DO KNEES CLICK
WHY AREN'T THERE E GRADES
WHY IS ISOLATION BAD
WHY DO BOYS LIKE ME
WHY DON'T BOYS LIKE ME
WHY IS THERE ALWAYS A JAVA UPDATE
WHY ARE THERE RED DOTS ON MY THIGHS
WHY IS LYING GOOD

Simpson's paradox

Edward Hugh Simpson

(10 December 1922 – 5 February 2019)

was a British codebreaker, statistician and civil servant.



"The Interpretation of Interaction in Contingency Tables", Journal of the Royal Statistical Society, 1951

Is it possible for one doctor to have a higher success rate than another doctor in every type of treatment he performs but to have a lower overall success rate across all treatment types?



Dr. Hibbert



Dr. Nick

Simpson's Paradox

	Hibbert heart bandaid	Nick heart bandaid
Success	70	2
Failure	20	8
	10	81
	0	9

Dr. Hibbert: success rate = 80%

Dr. Nick: success rate = 83%

Simpson's paradox

Edward Hugh Simpson

(10 December 1922 – 5 February 2019)

was a British codebreaker, statistician and civil servant.



"The Interpretation of Interaction in Contingency Tables", Journal of the Royal Statistical Society, 1951

Is it possible for one doctor to have a higher success rate than another doctor in every type of treatment he performs but to have a lower overall success rate across all treatment types?



Dr. Hibbert



Dr. Nick

Simpson's Paradox

	Hibbert heart bandaid	Nick heart bandaid
Success	70	2
Failure	20	8
	10	81
	0	9

Dr. Hibbert: success rate = 80%

Dr. Nick: success rate = 83%

Simpson's paradox might explain altruism

- Darwinian evolution has a problem with altruism
- “Selfish genes” do not care about others
- J. B. S. Haldane, (1892-1964)
British geneticist, evolutionary biologist
- When asked if he would give his life to save a drowning brother answered: “No, but I would to save two brothers or eight cousins”
- Altruism in some insect colonies like ants is because they are all genetically similar.



Altruism in bacteria

- Bacteria live in communities in close proximity to each other
- Individual bugs **spend significant resources** to produce **extracellular molecules**, excrete them outside of the cell to **share with others. That slows their growth**
 - Examples: extracellular enzymes, biofilm components, antimicrobial and anti-immune agents
- **Cheaters have faster growth rate**
 - They can take over by not producing any shared molecules
- **Evolutionary paradox: how bacteria can be altruistic?**

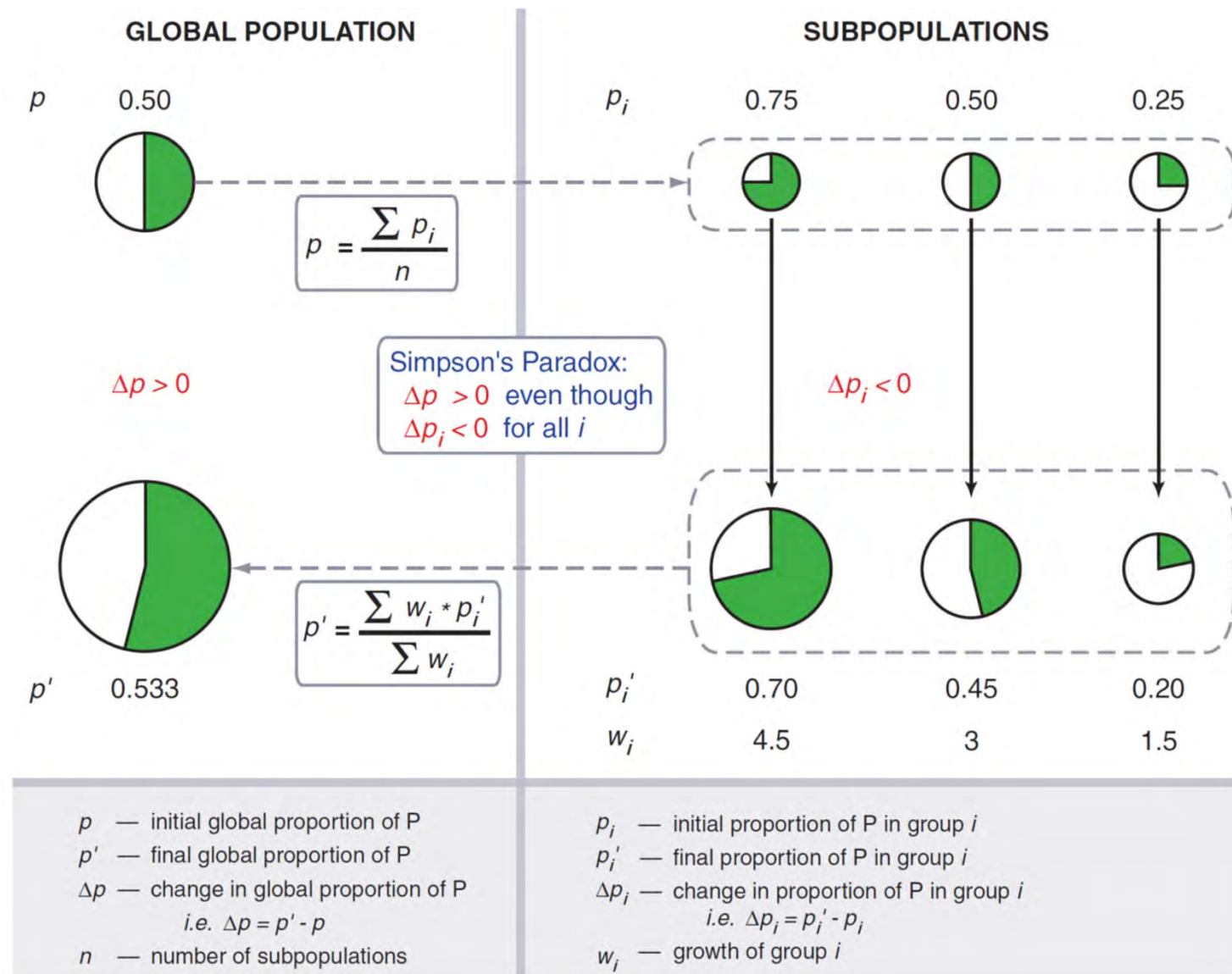


Simpson's Paradox in a Synthetic Microbial System

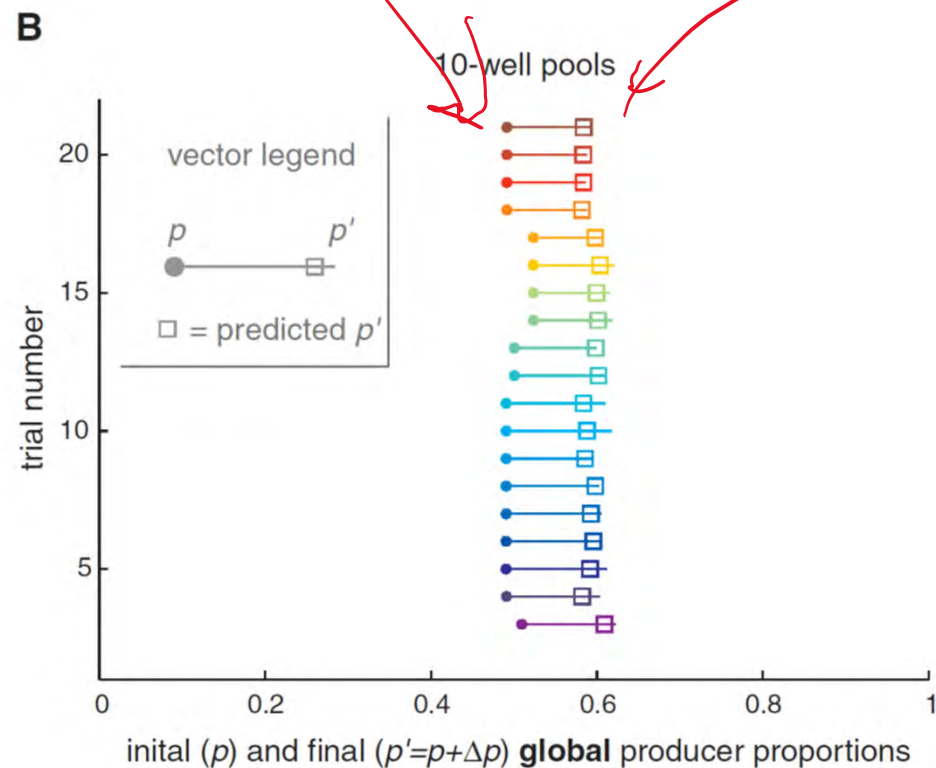
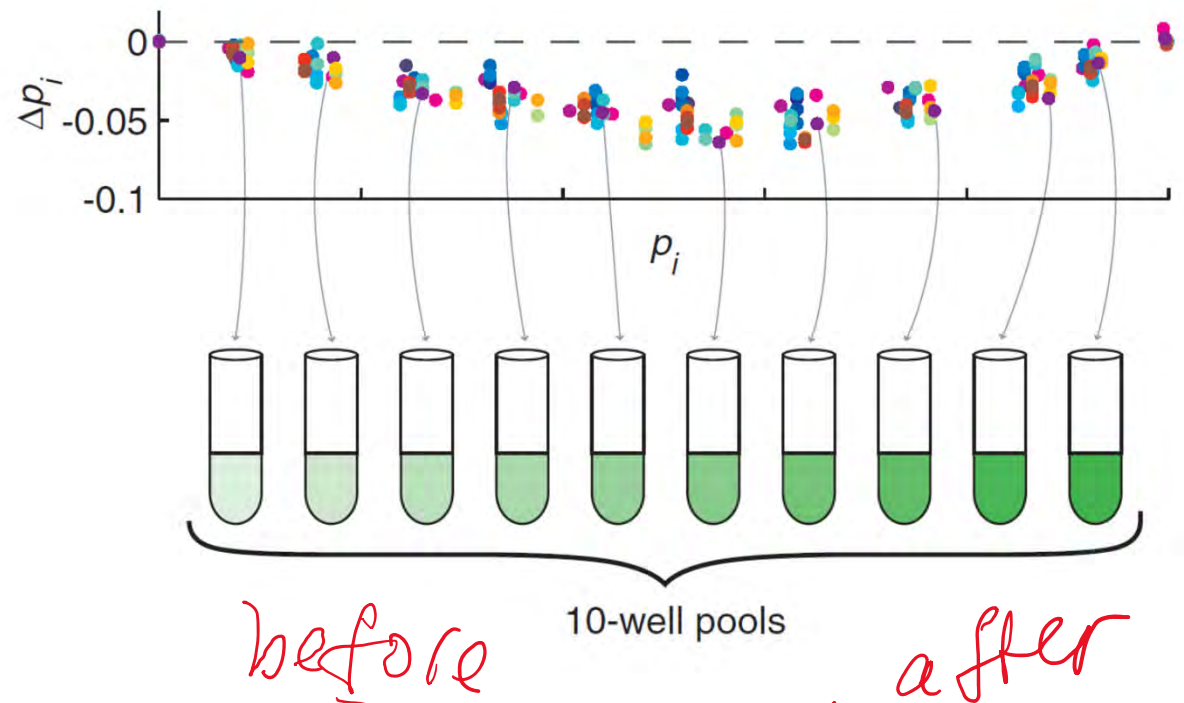
John S. Chuang,* Olivier Rivoire, Stanislas Leibler

The maintenance of “public” or “common good” producers is a major question in the evolution of cooperation. Because nonproducers benefit from the shared resource without bearing its cost of production, they may proliferate faster than producers. We established a synthetic microbial system consisting of two *Escherichia coli* strains of common-good producers and nonproducers. Depending on the population structure, which was varied by forming groups with different initial compositions, an apparently paradoxical situation could be attained in which nonproducers grew faster within each group, yet producers increased overall. We show that a simple way to generate the variance required for this effect is through stochastic fluctuations via population bottlenecks. The synthetic approach described here thus provides a way to study generic mechanisms of natural selection.

- The common good was a membrane-permeable Rhl autoinducer molecule rewired to activate antibiotic (chloramphenicol; Cm) resistance gene expression.



Fraction of altruists in
each of individual
test tubes dropped



Yet the overall fraction of
altruists in
all test tubes combined
increased

Credit: XKCD
comics

WHY ARE THERE SLAVES IN THE BIBLE

WHY DO TWINS HAVE DIFFERENT FINGERPRINTS
WHY ARE AMERICANS AFRAID OF DRAGONS

WHY IS HTTPS CROSSED OUT IN RED
WHY IS THERE A LINE THROUGH HTTPS
WHY IS THERE A RED LINE THROUGH HTTPS ON FACEBOOK
WHY IS HTTPS IMPORTANT

WHY AREN'T MY
ARMS GROWING



WHY ARE THERE WEEKS
WHY DO I FEEL DIZZY

QUESTIONS

FOUND IN GOOGLE AUTOCOMPLETE

WHY AREN'T ECONOMISTS RICH

WHY DO AMERICANS CALL IT SOCCER

WHY ARE MY EARS RINGING

WHY ARE THERE SO MANY AVENGERS

WHY ARE THE AVENGERS FIGHTING THE X MEN

WHY IS WOLVERINE NOT IN THE AVENGERS

WHY ARE THERE ANTS IN MY LAPTOP

WHY IS EARTH TILTED

WHY IS SPACE BLACK

WHY IS OUTER SPACE SO COLD

WHY ARE THERE PYRAMIDS ON THE MOON

WHY IS NASA SHUTTING DOWN

WHY ARE THERE MALE AND FEMALE BIKES

WHY ARE THERE TINY SPIDERS IN MY HOUSE

WHY DO SPIDERS COME INSIDE

WHY ARE THERE HUGE SPIDERS IN MY HOUSE

WHY ARE THERE LOTS OF SPIDERS IN MY HOUSE

WHY ARE THERE SPIDERS IN MY ROOM

WHY ARE THERE SO MANY SPIDERS IN MY ROOM

WHY DO SPIDER BITES ITCH

WHY IS DYING SO SCARY

WHY IS THERE NO GPS IN LAPTOPS

WHY DO KNEES CLICK

WHY AREN'T THERE E GRADES

WHY IS SEX
SO IMPORTANT



WHY ARE THERE
GHOSTS



WHY IS THERE AN OWL IN MY BACKYARD

WHY IS THERE AN OWL OUTSIDE MY WINDOW

WHY IS THERE AN OWL ON THE DOLLAR BILL

WHY DO OWLS ATTACK PEOPLE

WHY ARE AK 47s SO EXPENSIVE

WHY ARE THERE HELICOPTERS CIRCLING MY HOUSE

WHY ARE THERE GODS

WHY ARE THERE TWO SPOCKS

WHY IS MT VESUVIUS THERE

WHY DO THEY SAY T MINUS

WHY ARE THERE OBELISKS

WHY ARE WRESTLERS ALWAYS WET

WHY ARE OCEANS BECOMING MORE ACIDIC

WHY IS ARWEN DYING

WHY AREN'T MY QUAIL LAYING EGGS

WHY AREN'T MY QUAIL EGGS HATCHING

WHY AREN'T THERE ANY FOREIGN MILITARY BASES IN AMERICA

WHY ARE CIGARETTES LEGAL
WHY ARE THERE DUCKS IN MY POOL
WHY IS JESUS WHITE
WHY IS THERE LIQUID IN MY EAR
WHY DO Q TIPS FEEL GOOD
WHY DO GOOD PEOPLE DIE

WHY AREN'T
THERE GUNS IN
HARRY POTTER



WHY ARE ULTRASOUNDS IMPORTANT
WHY ARE ULTRASOUND MACHINES EXPENSIVE
WHY IS STEALING WRONG

WHY ARE DOGS AFRAID OF FIREWORKS
WHY IS THERE NO KING IN ENGLAND

WHY DO WHALES JUMP
WHY ARE WITCHES GREEN
WHY ARE THERE MIRRORS ABOVE BEDS
WHY DO I SAY UH

WHY IS SEA SALT BETTER

WHY ARE THERE TREES IN THE MIDDLE OF FIELDS

WHY IS THERE NOT A POKEMON MMO

WHY IS THERE LAUGHING IN TV SHOWS

WHY ARE THERE DOORS ON THE FREEWAY

WHY ARE THERE SO MANY SVCHOST.EXE RUNNING

WHY AREN'T THERE ANY COUNTRIES IN ANTARCTICA

WHY ARE THERE SCARY SOUNDS IN MINECRAFT

WHY IS THERE KICKING IN MY STOMACH

WHY ARE THERE TWO SLASHES AFTER HTTP

WHY ARE THERE CELEBRITIES

WHY DO SNAKES EXIST

WHY DO OYSTERS HAVE PEARLS

WHY ARE DUCKS CALLED DUCKS

WHY DO THEY CALL IT THE CLAP

WHY ARE KYLE AND CARTMAN FRIENDS

WHY IS THERE AN ARROW ON AANG'S HEAD

WHY ARE TEXT MESSAGES BLUE

WHY ARE THERE MUSTACHES ON CLOTHES

WHY ARE THERE MUSTACHES ON CARS

WHY ARE THERE MUSTACHES EVERYWHERE

WHY ARE THERE SO MANY BIRDS IN OHIO

WHY IS THERE SO MUCH RAIN IN OHIO

WHY IS OHIO WEATHER SO WEIRD

WHY ARE THERE BRIDESMAIDS

WHY DO DYING PEOPLE REACH UP

WHY AREN'T THERE VARICOSE ARTERIES

WHY ARE OLD KINGDOMS DIFFERENT

WHY ARE THERE
SQUIRRELS



WHY IS PROGRAMMING SO HARD

WHY IS THERE A 0 OHM RESISTOR

WHY DO AMERICANS HATE SOCCER

WHY DO RHYMES SOUND GOOD

WHY DO TREES DIE

WHY IS THERE NO SOUND ON CNN

WHY AREN'T POKEMON REAL

WHY AREN'T BULLETS SHARP

WHY DO DREAMS SEEM SO REAL

WHY AREN'T THERE DINOSAUR GHOSTS

WHY ARE THERE FEMALE MR NIMES

WHY IS GPS FREE

WHY IS LIVING GOOD

Let's check the theory by playing the game

Go to

<https://dacalderon.shinyapps.io/montyhall/>

- Tables 1,3,5 will play “switch the door” strategy
- Tables 2,4,6 will play “same door” strategy
- Play at least 30 rounds (more is better)
- In the end we will **add up the numbers from all tables**

Let's check with more random experiments

- Stats=??;
- %set Stats large...
- switch_count=0; noswitch_count=0; %set 0 at the beginning
- for n = 1:Stats
- a = randperm(3); %Monty places two goats and the car at random
- %a(1) -goat, a(2) -goat, a(3) - car
- i= floor(3.*rand)+1; %you select the door!
- % SWITCH STRATEGY
- if(i == a(1)) switch_count=switch_count+??; %a(2)-opened, switch to a(3), car!
- elseif (i == a(2)) switch_count = switch_count + ??;%a(1) opened, switch to a(3), car!
- else switch_count = switch_count + ??; %a(1)/a(2) opened, switch to a(2)/a(1), no car :-(
- end
- % NO SWITCH STRATEGY
- if(i == a(1)) noswitch_count = noswitch_count + ??; %a(2)-opened, no car :-(
- elseif (i==a(2)) noswitch_count = noswitch_count + ?? %a(1)-opened, no car :-(
- else noswitch_count = noswitch_count + ??; %a(1) or a(2)-opened, car!
- endend;
- disp('probability to win a car if switched doors=');
- disp(num2str(switch_count./??)); %# of cars with switching
- disp('probability to win a car if did not switch doors=');
- disp(num2str(noswitch_count./??)); %# of cars w/o switching



Discrete Probability Distributions

Random Variables

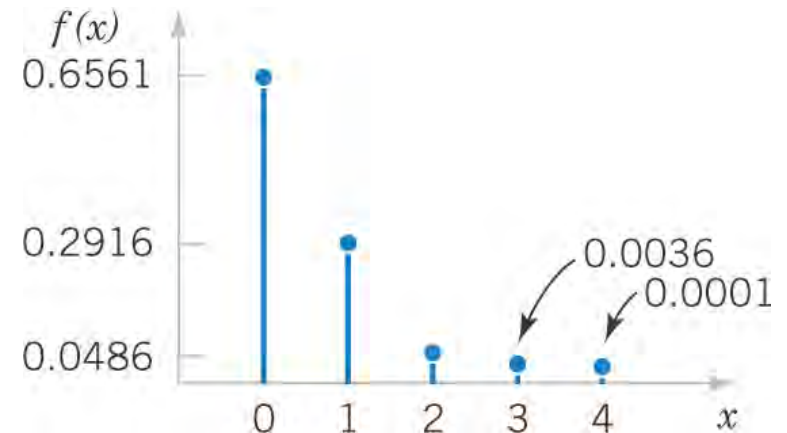
- A variable that associates a number with the outcome of a **random experiment** is called a **random variable**.
- Notation: **random variable** is denoted by an uppercase letter, such as ***X***. After the experiment is conducted, the **measured value** is denoted by a **lowercase letter**, such as ***x***. Both *X* and *x* are shown in italics, e.g., ***P(X=x)***.

Continuous & Discrete Random Variables

- A **discrete random variable** is usually integer number
 - N - the number of p53 proteins in a cell
 - D - the number of nucleotides different between two sequences
- A **continuous random variable** is a real number
 - $C=N/V$ – the concentration of p53 protein in a cell of volume V
 - Percentage $(D/L)*100\%$ of different nucleotides in protein sequences of different lengths L
(depending on the set of L 's may be discrete but dense)

Probability Mass Function (PMF)

- I want to **compare all 4-mers** in a pair of human genomes
- **X – random variable:** the number of nucleotide differences in a given 4-mer
- **Probability Mass Function:** $f(x)$ or $P(X=x)$ – the probability that the # of SNPs is **exactly equal to x**



Probability Mass Function for the # of mismatches in 4-mers

$P(X=0) =$	0.6561
$P(X=1) =$	0.2916
$P(X=2) =$	0.0486
$P(X=3) =$	0.0036
$P(X=4) =$	0.0001
$\sum_x P(X=x) =$	1.0000

Cumulative Distribution Function (CDF)

	x	$P(X=x)$	$P(X \leq x)$	$P(X > x)$
	-1	0.0000	0.0000	1.0000
	0	0.6561	0.6561	0.3439
	1	0.2916	0.9477	0.0523
	2	0.0486	0.9963	0.0037
	3	0.0036	0.9999	0.0001
	4	0.0001	1.0000	0.0000

Cumulative Distribution Function CDF: $F(x) = P(X \leq x)$

Example:

$$F(3) = P(X \leq 3) = P(X=0) + P(X=1) + P(X=2) + P(X=3) = 0.9999$$

Complementary Cumulative Distribution Function
(tail distribution) or CCDF: $F_{>}(x) = P(X > x)$

$$\text{Example: } F_{>}(0) = P(X > 0) = 1 - P(X \leq 0) = 1 - 0.6561 = 0.3439$$

Mean or Expected Value of X

The **mean** or **expected value** of the discrete random variable X, denoted as μ or $E(X)$, is

$$\mu = E(X) = \sum_x x \cdot P(X = x) = \sum_x x \cdot f(x)$$

- **The mean** = the weighted average of all possible values of X. It represents its “center of mass”
- The **mean** may, or may not, be an **allowed value of X**
- It is also called the **arithmetic mean** (to distinguish from e.g. the **geometric mean** discussed later)
- **Mean** may be infinite if X any integer and tail $P(X=x) > c/x^2$

Outcomes of 6 random experiments

0, 1, 0, 0, 2, 1

$$\text{Mean} = \frac{0 + 1 + 0 + 0 + 2 + 1}{6} =$$

$$= \frac{3 \times 0 + 2 \times 1 + 1 \times 2}{6} =$$

$$= 0 \times \frac{3}{6} + 1 \times \frac{2}{6} + 2 \times \frac{1}{6} = \sum_{x=0}^2 x P(X=x)$$

- $E[X] = \sum_x x \cdot P(X=x)$

- $E[X^2] = \sum_x x^2 \cdot P(X=x)$

- $E[a \cdot X + b \cdot X^2] = \sum (ax + bx^2) \cdot$
 $\cdot P(X=x) = a \cdot \sum x P(X=x) +$
 $+ b \sum x^2 P(X=x)$

- $E[e^X] = \sum e^x P(X=x)$

Variance $V(X)$: Square
of a typical deviation from
the mean $\mu = E(X)$

$V(X) = \sigma^2$, where σ is called
standard deviation

$$\begin{aligned}\sigma^2 &= V(X) = E((X - \mu)^2) = \\&= E(X^2 - 2\mu X + \mu^2) = E(X^2) - \\&- 2\mu E(X) + \mu^2 = E(X^2) - 2\mu^2 + \mu^2 = \\&= E(X^2) - \mu^2 = E(X^2) - (E(X))^2\end{aligned}$$

Variance of a Random Variable

If X is a discrete random variable with probability mass function $f(x)$,

$$E[h(X)] = \sum_x h(x) \cdot P(X = x) = \sum_x h(x) f(x) \quad (3-4)$$

If $h(x) = (X - \mu)^2$, then its expectation, $V(x)$, is the **variance of X** .

$\sigma = \sqrt{V(x)}$, is called **standard deviation of X**

$\sigma^2 = V(X) = \sum_x (x - \mu)^2 f(x)$ is the **definitional** formula

$$= \sum_x (x^2 - 2\mu x + \mu^2) f(x)$$

$$= \sum_x x^2 f(x) - 2\mu \sum_x x f(x) + \mu^2 \sum_x f(x)$$

$$= \sum_x x^2 f(x) - 2\mu^2 + \mu^2$$

$$= \sum_x x^2 f(x) - \mu^2 \text{ is the } \mathbf{computational} \text{ formula}$$

Variance can be infinite
if X can be any integer
and tail of $P(X=x) \geq c/x^3$

Skewness of a random variable

- Want to quantify **how asymmetric** is the **distribution around the mean?**
- Need any **odd moment**: $E[(X-\mu)^{2n+1}]$
- **Cannot** do it with the **first moment**: $E[X-\mu]=0$
- Normalized 3-rd moment is **skewness**: $\gamma_1 = E[(X-\mu)^3]/\sigma^3$
- Skewness **can be infinite** if X takes unbounded integer values and tail $P(X=x) \geq c/x^4$

Geometric mean of a random variable

- Useful for **very broad distributions** (many orders of magnitude)?
- Mean may be dominated by **very unlikely** but **very large events**. Think of a **lottery**
- **Exponent of the mean of $\log X$:**
Geometric mean = $\exp(E[\log X])$
- Geometric mean usually **is not infinite**

Summary: Parameters of a Probability Distribution

- **Probability Mass Function (PMF):** $f(x) = \text{Prob}(X=x)$
- **Cumulative Distribution Function (CDF):** $F(x) = \text{Prob}(X \leq x)$
- **Complementary Cumulative Distribution Function (CCDF):**
 $F_{>}(x) = \text{Prob}(X > x)$
- The **mean**, $\mu = E[X]$, is a measure of the **center of mass of a random variable**
- The **variance**, $V(X) = E[(X - \mu)^2]$, is a measure of the **dispersion** of a random variable **around its mean**
- The **standard deviation**, $\sigma = [V(X)]^{1/2}$, is another measure of the **dispersion** around mean. Has the same units as X
- The **skewness**, $\gamma_1 = E[(X - \mu)^3 / \sigma^3]$, a measure of asymmetry around mean
- The **geometric mean**, $\exp(E[\log X])$ is useful for very broad distributions

Skewness of a random variable

- Want to quantify **how asymmetric** is the **distribution around the mean?**
- Need any **odd moment**: $E[(X-\mu)^{2n+1}]$
- **Cannot** do it with the **first moment**: $E[X-\mu]=0$
- Normalized 3-rd moment is **skewness**: $\gamma_1 = E[(X-\mu)^3/\sigma^3]$
- Skewness **can be infinite** if X takes unbounded positive integer values and the tail $P(X=x) \geq c/x^4$ for large x

Geometric mean of a random variable

- Useful for **very broad distributions** (many orders of magnitude)?
- Mean may be dominated by **very unlikely** but **very large events**. Think of a **lottery**
- **Exponent of the mean of $\log X$:**
Geometric mean = $\exp(E[\log X])$
- Geometric mean usually **is not infinite**

Summary: Parameters of a Probability Distribution

- **Probability Mass Function (PMF):** $f(x)=\text{Prob}(X=x)$
- **Cumulative Distribution Function (CDF):** $F(x)=\text{Prob}(X\leq x)$
- **Complementary Cumulative Distribution Function (CCDF):** $F_{>}(x)=\text{Prob}(X>x)$
- The **mean**, $\mu=E[X]$, is a measure of the **center of mass of a random variable**
- The **variance**, $V(X)=E[(X-\mu)^2]$, is a measure of the **dispersion** of a random variable **around its mean**
- The **standard deviation**, $\sigma=[V(X)]^{1/2}$, is another measure of the **dispersion** around mean. Has the same units as X
- The **skewness**, $\gamma_1=E[(X-\mu)^3/\sigma^3]$, a measure of asymmetry around mean
- The **geometric mean**, $\exp(E[\log X])$ is useful for very broad distributions

A gallery of useful discrete probability distributions

Discrete Uniform Distribution

- Simplest discrete distribution.
- The random variable X assumes only a finite number of values, each with equal probability.
- A random variable X has a discrete uniform distribution if each of the n values in its range, say x_1, x_2, \dots, x_n , has equal probability.

$$f(x_i) = 1/n$$

Uniform Distribution of Consecutive Integers

- Let X be a discrete uniform random variable all integers from a to b (inclusive). There are $b - a + 1$ integers. Therefore each one gets:

$$f(x) = 1/(b-a+1)$$

- Its measures are:

$$\mu = E(x) = (b+a)/2$$

$$\sigma^2 = V(x) = [(b-a+1)^2-1]/12$$

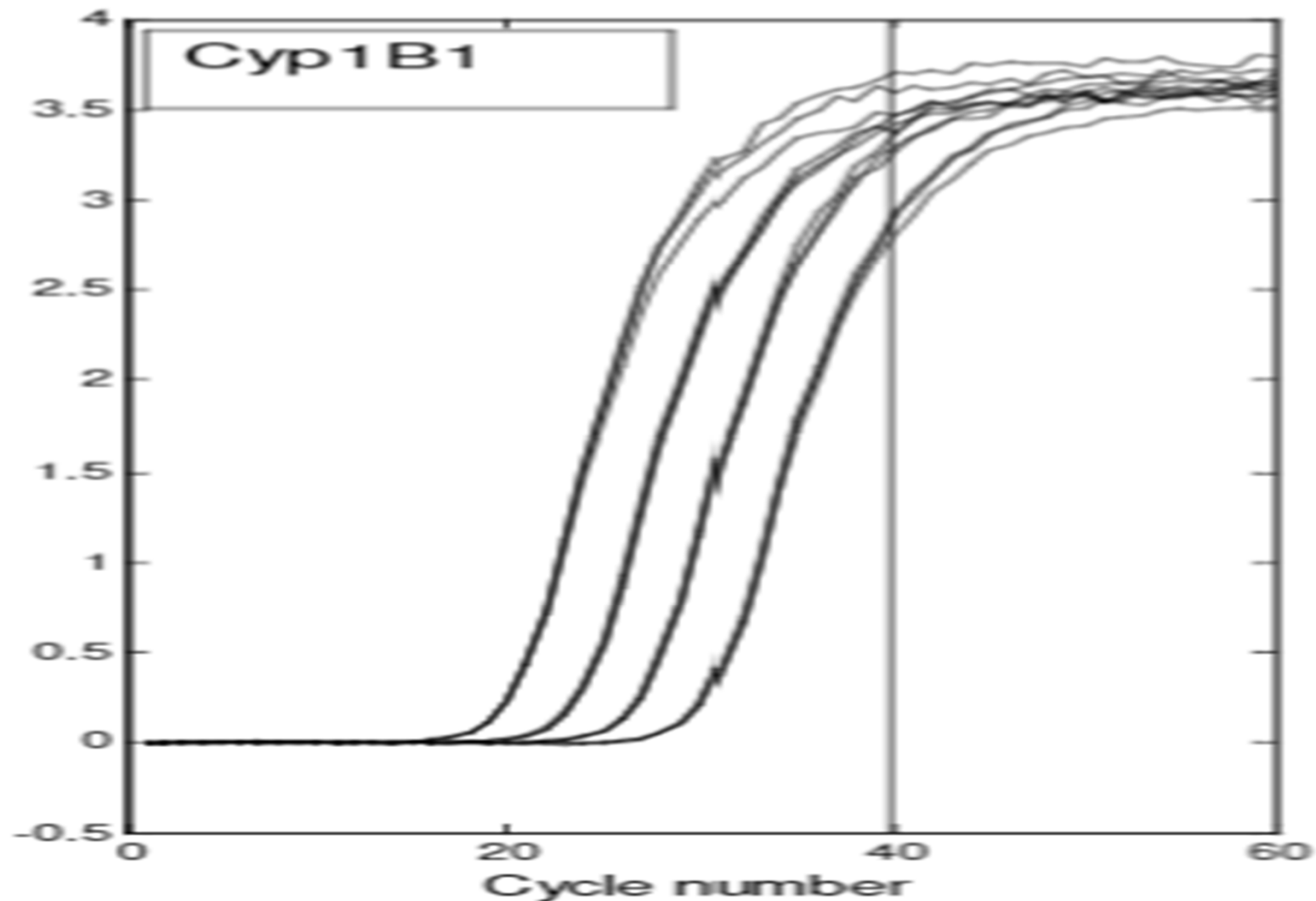
Note that the mean is the midpoint of a & b .

An example of the uniform
distribution

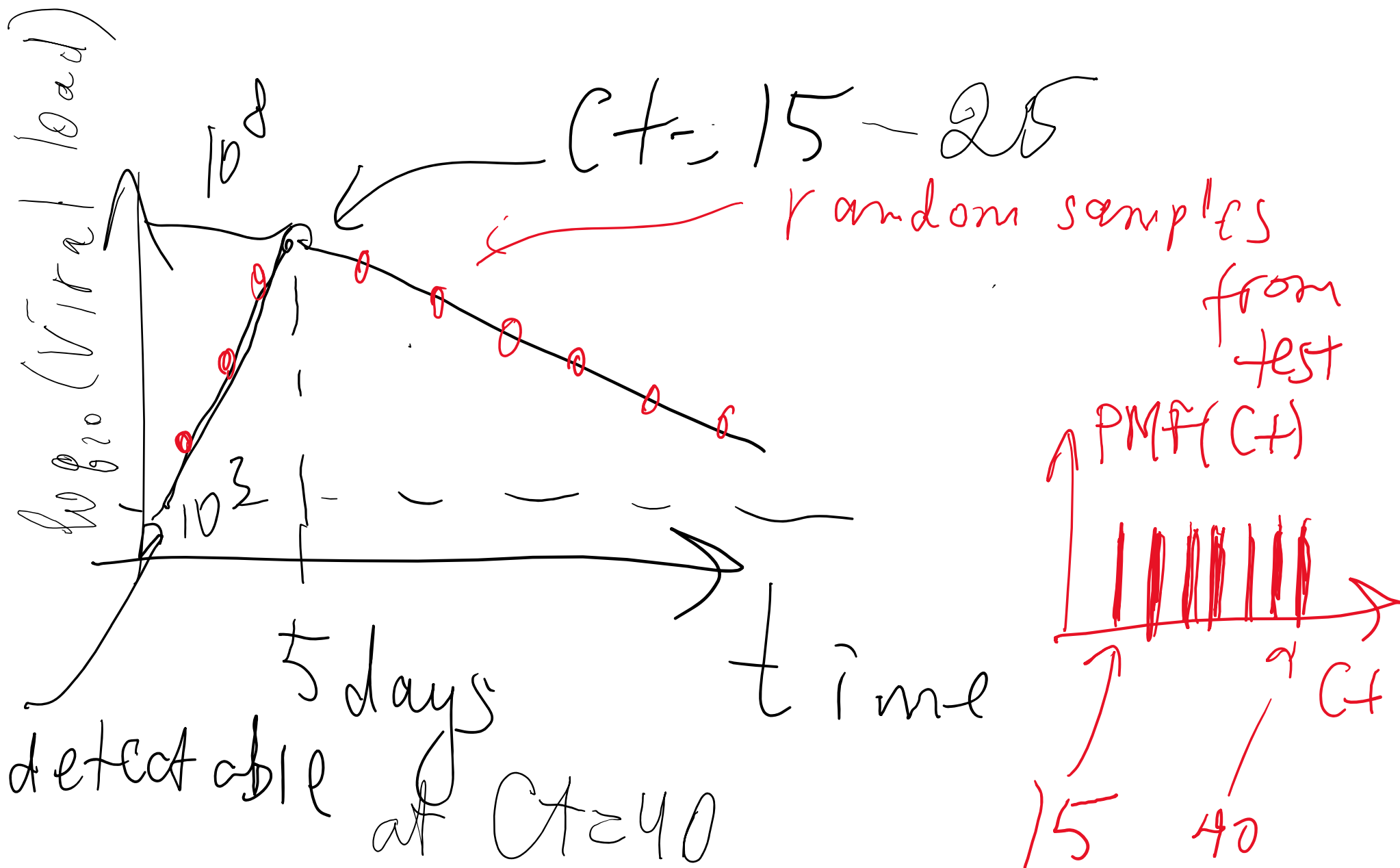
Cycle threshold (Ct) value in
COVID-19 infection

What is the Ct value of a PCR test?

$Ct = \text{const} - \log_2(\text{viral DNA concentration})$



Why Ct distribution should it be uniform?



Examples of uniform distribution: Ct value of PCR test of a virus

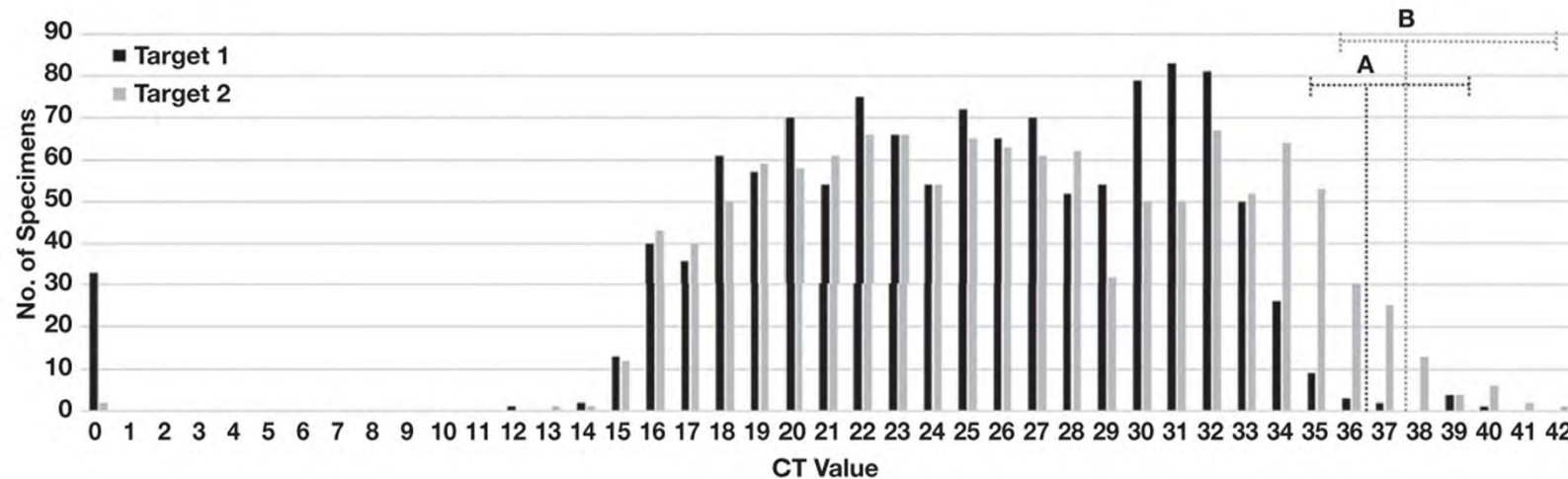


Figure 3 Distribution of cycle threshold (CT) values. The total number of specimens with indicated CT values for Target 1 and 2 are plotted. The estimated limit of detection for (A) Target 1 and (B) Target 2 are indicated by vertical dotted lines. Horizontal dotted lines encompass specimens with CT values less than 3x the LoD for which sensitivity of detection may be less than 100%. This included 19/1,180 (1.6%) reported CT values for Target 1 and 81/1,211 (6.7%) reported CT values for Target 2. Specimens with Target 1 or 2 reported as “not detected” are denoted as a CT value of “0.”

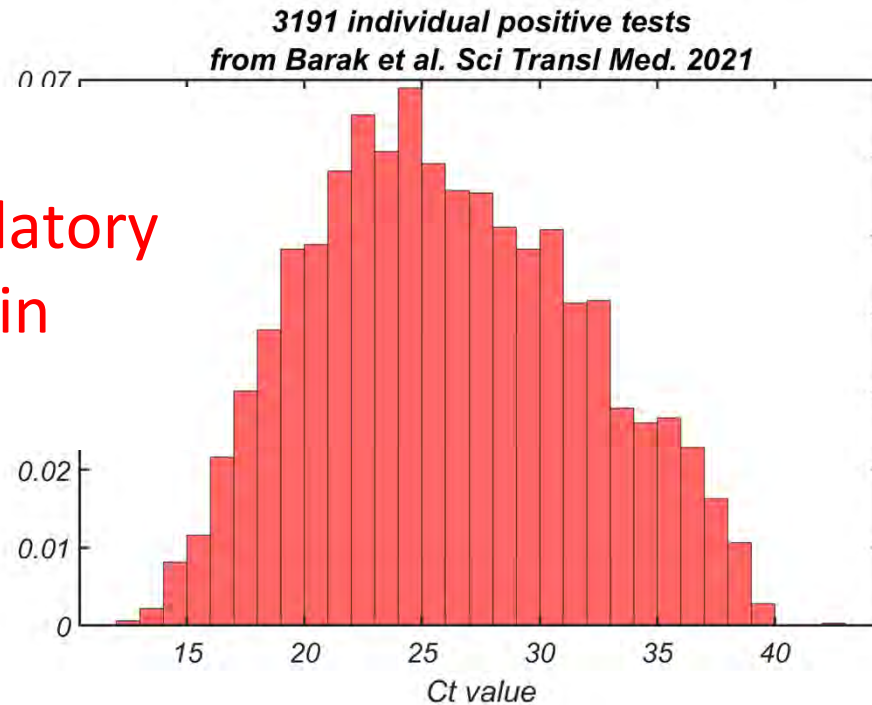
Distribution of SARS-CoV-2 PCR Cycle Threshold Values Provide Practical Insight Into Overall and Target-Specific Sensitivity Among Symptomatic Patients

Blake W Buchan, PhD, Jessica S Hoff, PhD, Cameron G Gmehlin, Adriana Perez, Matthew L Faron, PhD, L Silvia Munoz-Price, MD, PhD, Nathan A Ledebor, PhD *American Journal of Clinical Pathology*, Volume 154, Issue 4, 1 October 2020,

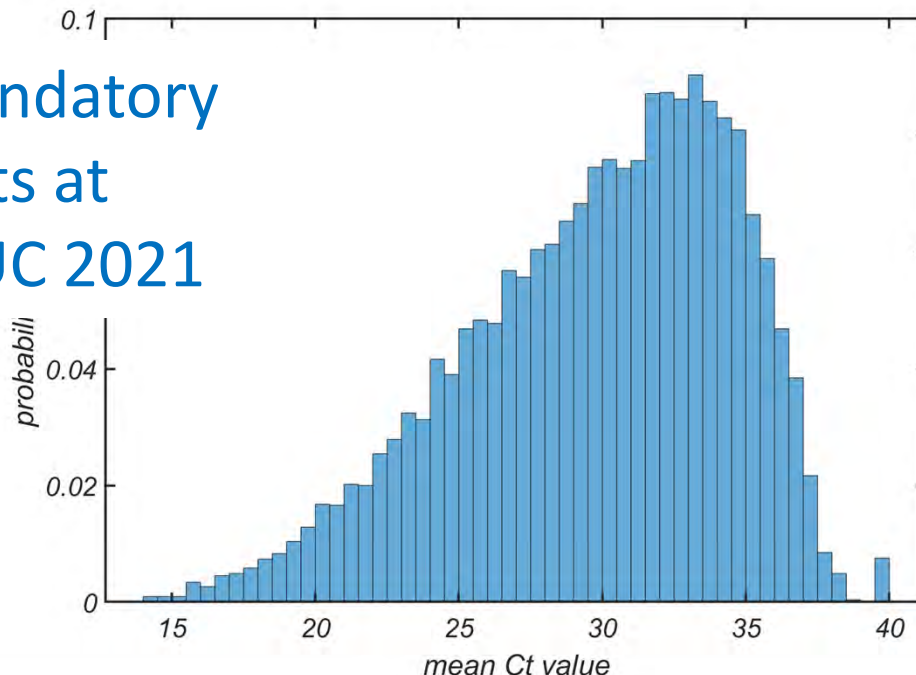
<https://academic.oup.com/ajcp/article/154/4/479/5873820>

Why should we care?

Non-mandatory tests in Israel



Mandatory tests at UIUC 2021



- High Ct value means we identified the infected individual early, hopefully before transmission to others
- When testing is mandatory, and people are tested frequently – Ct value is skewed towards high values

Matlab exercise: Uniform distribution

- Generate a **sample of size 100,000** for uniform random variable X taking values $1, 2, 3, \dots, 10$
- Plot the approximation to the probability mass function based on this sample
- Calculate mean and variance of this sample and compare it to infinite sample predictions:
 $E[X] = (a+b)/2$ and $V[X] = ((a-b+1)^2 - 1)/12$

Matlab template: Uniform distribution

- `b=10; a=1; % b= upper bound; a= lower bound (inclusive)'`
- `Stats=100000; % sample size to generate`
- `r1=rand(Stats,1);`
- `r2=floor(??*r1)+??;`
- `mean(r2)`
- `var(r2)`
- `std(r2)`
- `[hy,hx]=hist(r2, 1:10); % hist generates histogram in bins 1,2,3...,10`
- `% hy - number of counts in each bin; hx - coordinates of bins`
- `p_f=hy./??; % normalize counts to add up to 1`
- `figure; plot(??,p_f, 'ko-'); ylim([0, max(p_f)+0.01]); % plot the PMF`