

Mutational Burden of Tumors

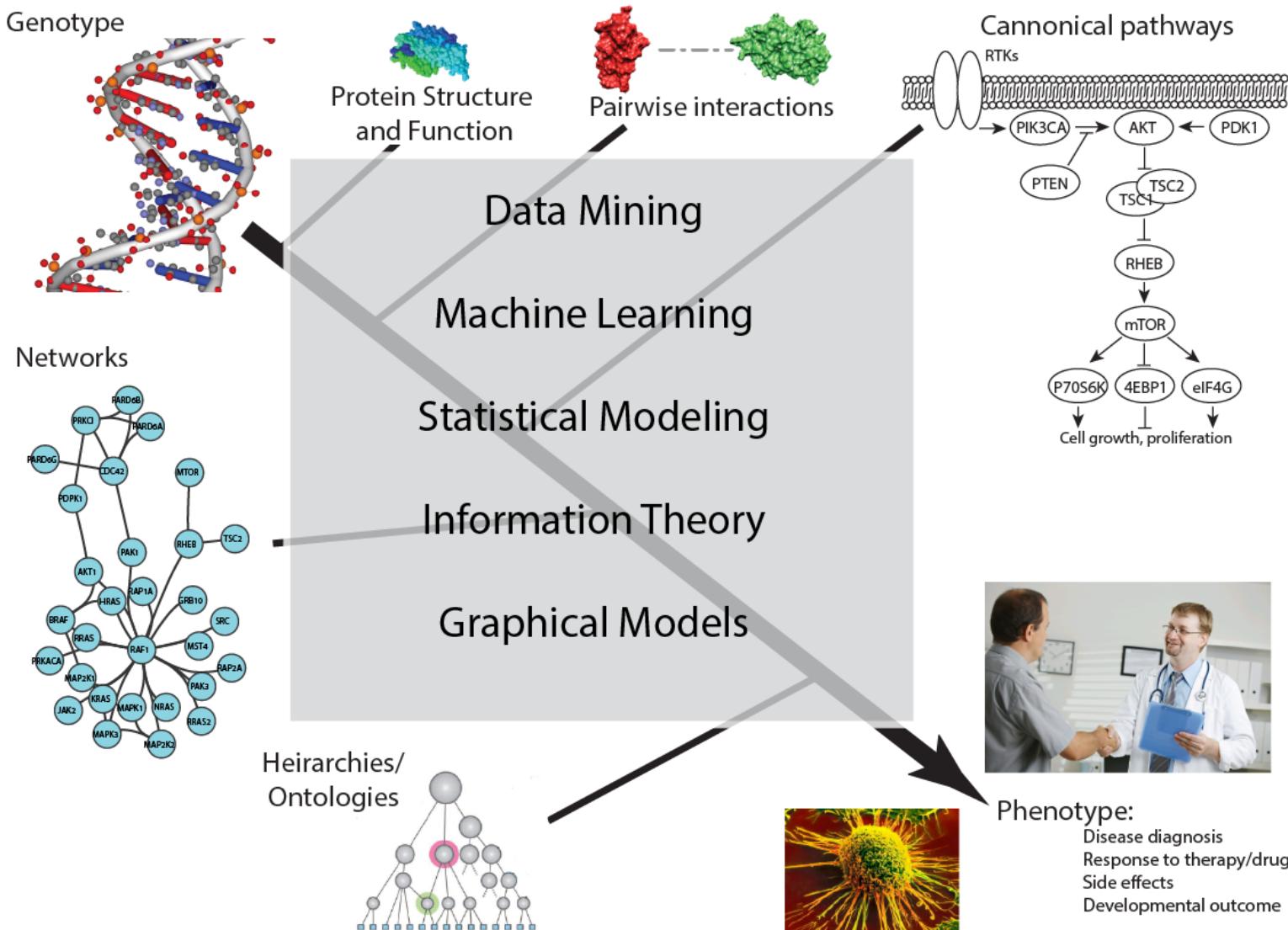
May 28, 2018

Hannah Carter, PhD

University of California, San Diego

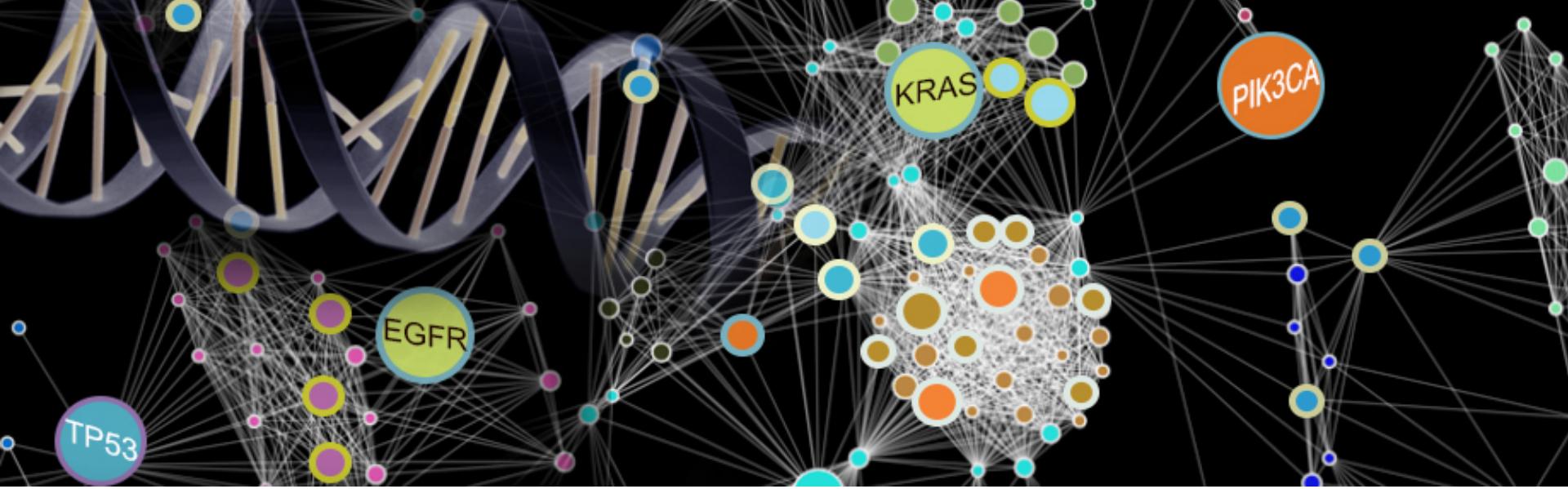
Department of Medicine

Studying mutations at multiple scales



Outline

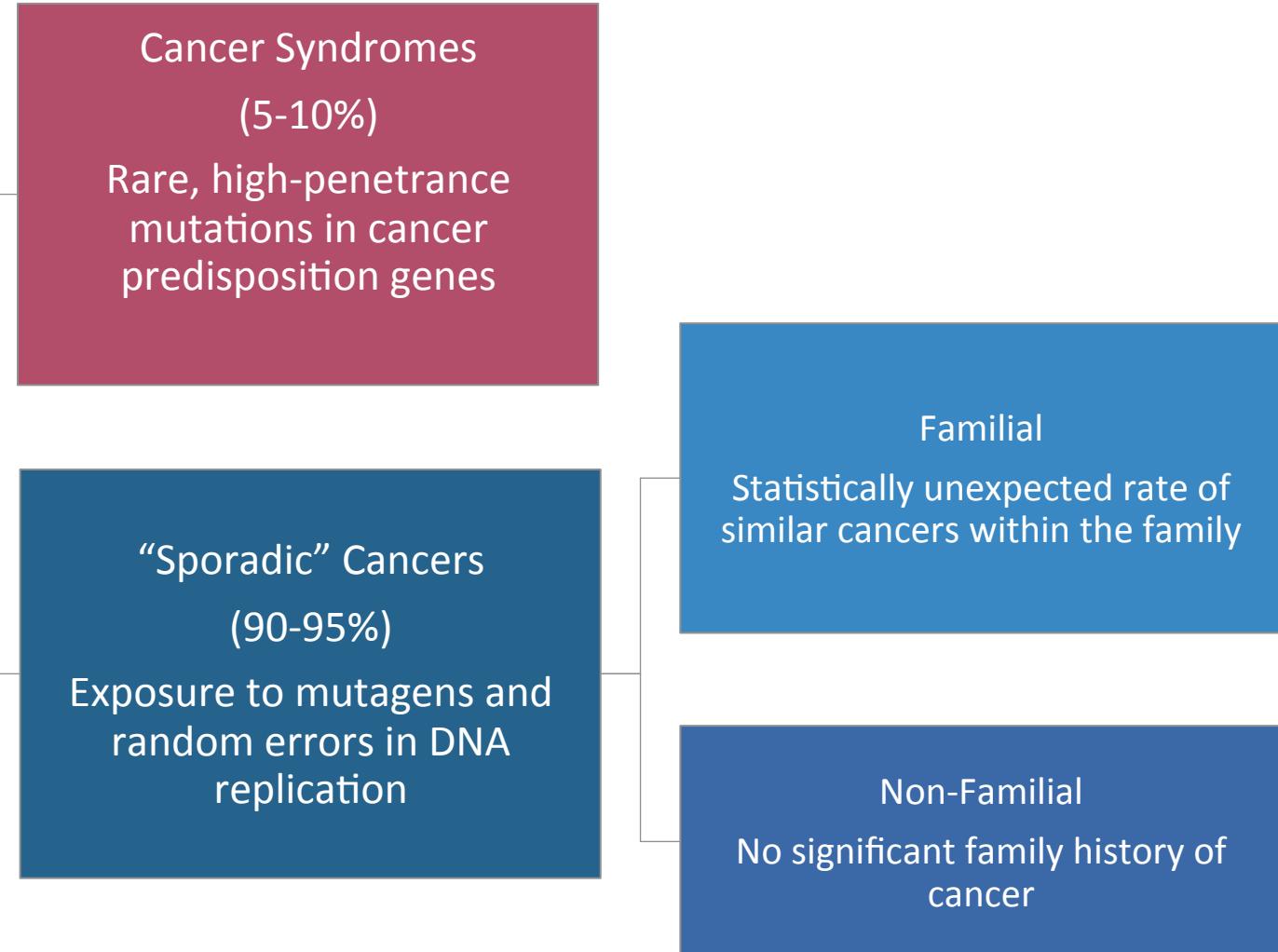
- Cancer as a genetic disease
- Mutation burden in tumors
- Identifying driver mutations in cancer
- Tumor heterogeneity
- Example: predicting driver mutations



Cancer as a genetic disease

Cancer

Cancer is a genetic disease



Familial cancer risk

Updated estimates of cancer heritability from a study of > 80,000 twins

Heritability of cancer overall was 33% (95% CI, 30%-37%)

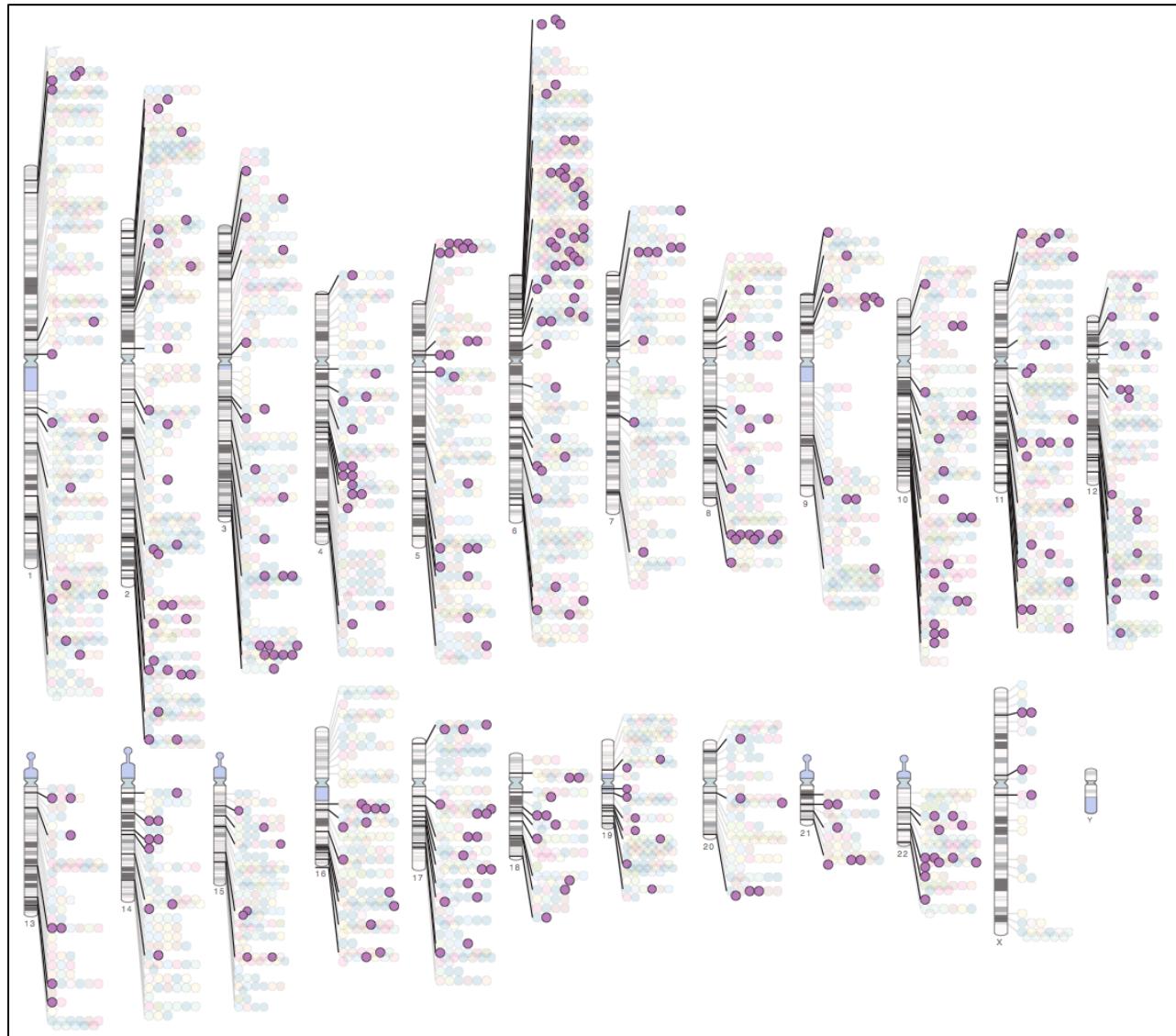
Cancer	Heritability	95% CI
Melanoma	58%	43-73%
Prostate	57%	51-63%
NonMelanoma Skin	43%	26-59%
Ovarian	39%	23-55%
Kidney	38%	21-55%
Breast	31%	11-51%
Corpus Uteri	27%	11-43%

Heritability: proportion of variance in cancer risk due to inter-individual genetic differences

Common cancer risk variants: Genome-Wide Association Studies

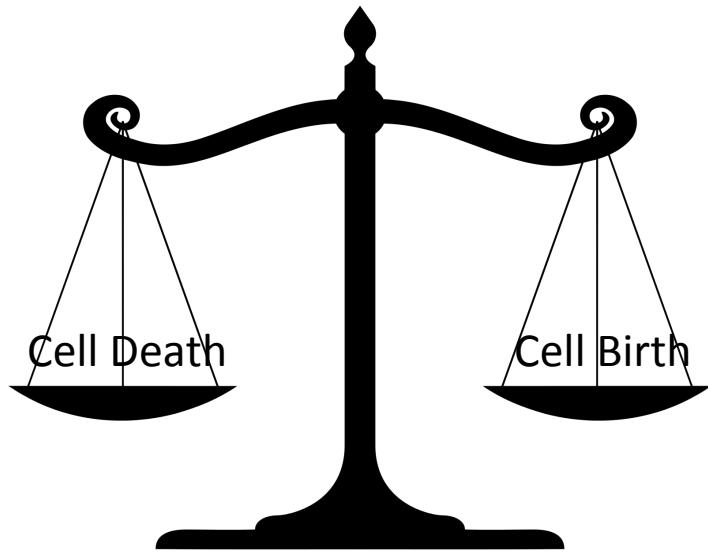
NHGRI GWAS Catalog:

- 165⁺ cancer studies
- Over 500 risk SNPs
- Population specific

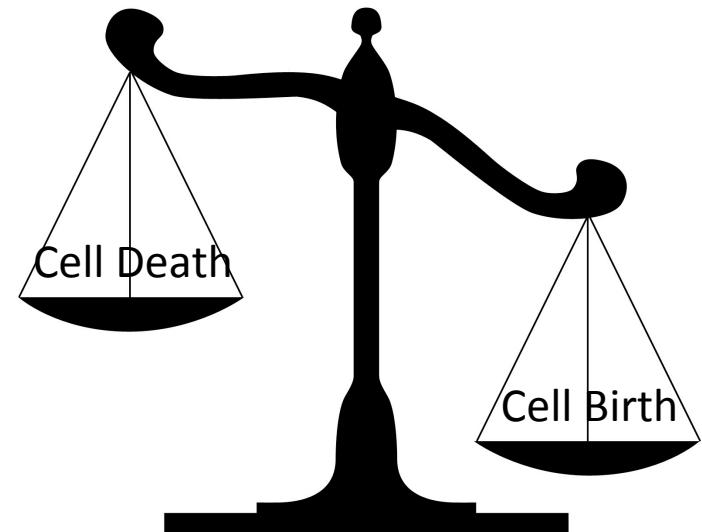


Somatic genetic alterations in cancer

Normal Tissues



Tumors

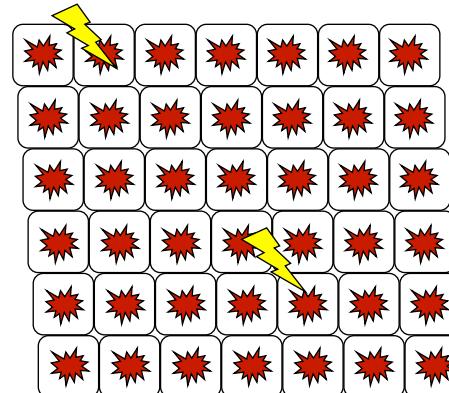
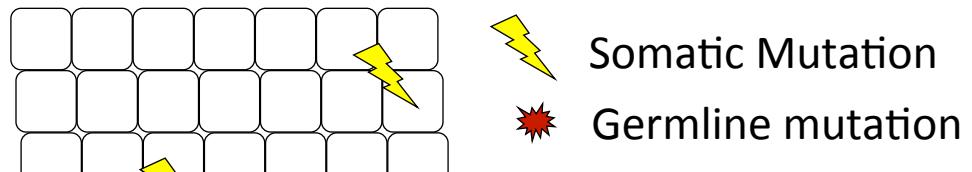


Accumulating genetic and epigenetic changes result in cancer by altering the balance between cell birth and cell death.

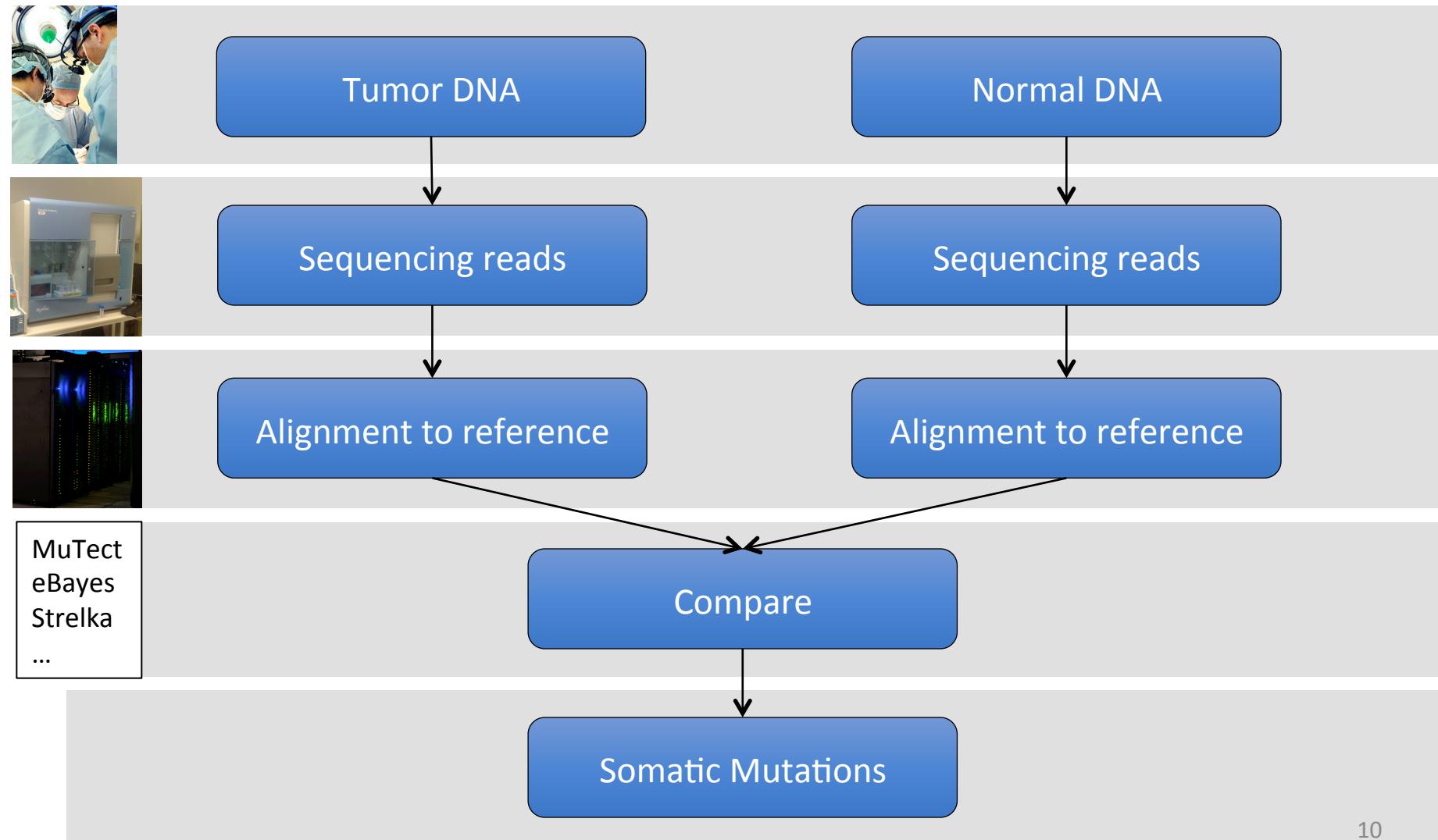
Concept: Germline-somatic interactions

Cancer results from multiple mutations co-occurring in the same cell

- Sporadic cancers
 - Somatic mutations only
- Genetic background
 - Germline variants could modify the fitness effects of somatic mutations



Cataloging mutations in tumors by DNA sequencing



Types of cancer mutations

Small-Scale Mutations

Point mutation	Change in a single nucleotide
Silent mutation	Amino acid sequence is not changed
Missense mutation	Amino acid sequence is changed
Nonsense mutation	Amino acid sequence is changed to a stop codon, therefore truncating the protein
Insertion	Addition of one or more extra nucleotides
Deletion	Removal of one or more nucleotides

Large-Scale Mutations

Amplification	Multiple copies of chromosomal region
Deletion	Loss of chromosomal region
Translocation	Interchange of regions from different chromosomes
Inversion	Reversal of the orientation of a chromosomal region
Loss of heterozygosity	Loss of one allele

Most tumor sequencing data is whole exome

Most analyses still focus predominantly on point mutations

Mutation burden = number of mutations per sequenced MB

Coding Base Substitutions

- Genetic Code

- 3 bases = 1AA
- “Degeneracy”
 - 20 AAs, 64 possible 3 base combinations
 - Not every base substitution will result in an AA change
 - 3rd position is least likely to result in a change

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

Coding Base Substitutions

Synonymous (or silent) Vs Non-synonymous

A) Silent

Wild Type — ACA — CAC — GAG — CCC —
Thr His Glu Pro

Mutant — ACA — CAC — GAA — CCC —
Thr His Glu Pro

C) Nonsense

— CAC — GAG — CCC — CTC —
His Glu Pro Leu
— CAC — TAG — CCC — CTC —
His STOP

B) Missense

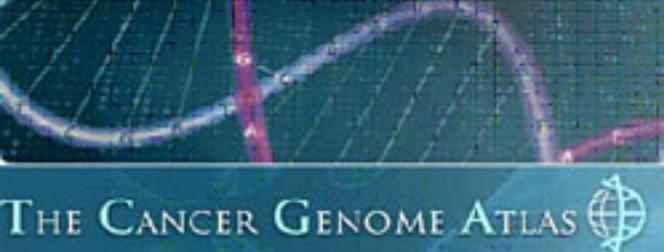
Wild Type — ACA — CAC — GAG — CCC —
Thr His Glu Pro

Mutant — ACA — CAC — GAC — CCC —
Thr His Asp Pro

D) Nonstop

— CCC — TAG — AAG — AGA —
Pro STOP

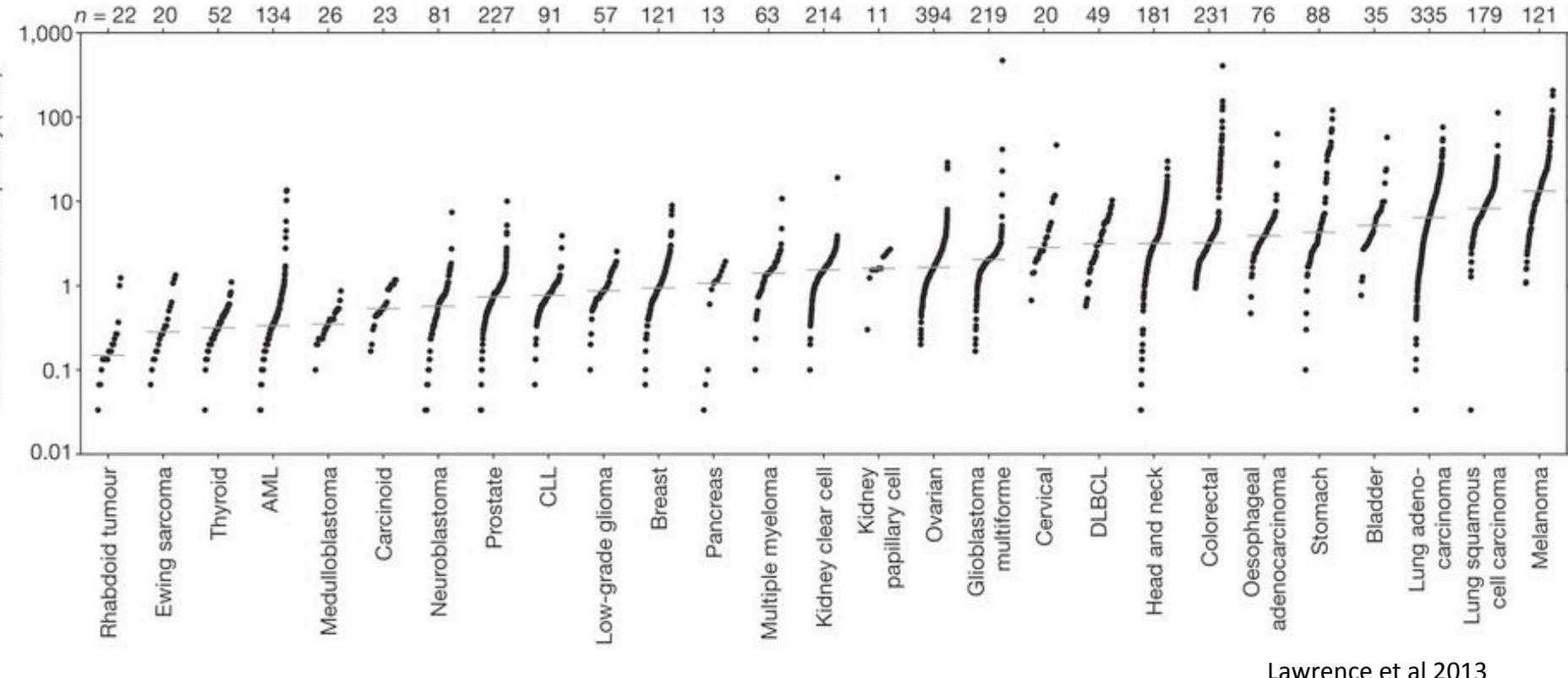
— CCC — TAC — AAG — AGA —
Pro Tyr Lys Arg



Somatic mutations uncovered by large-scale tumor exome sequencing

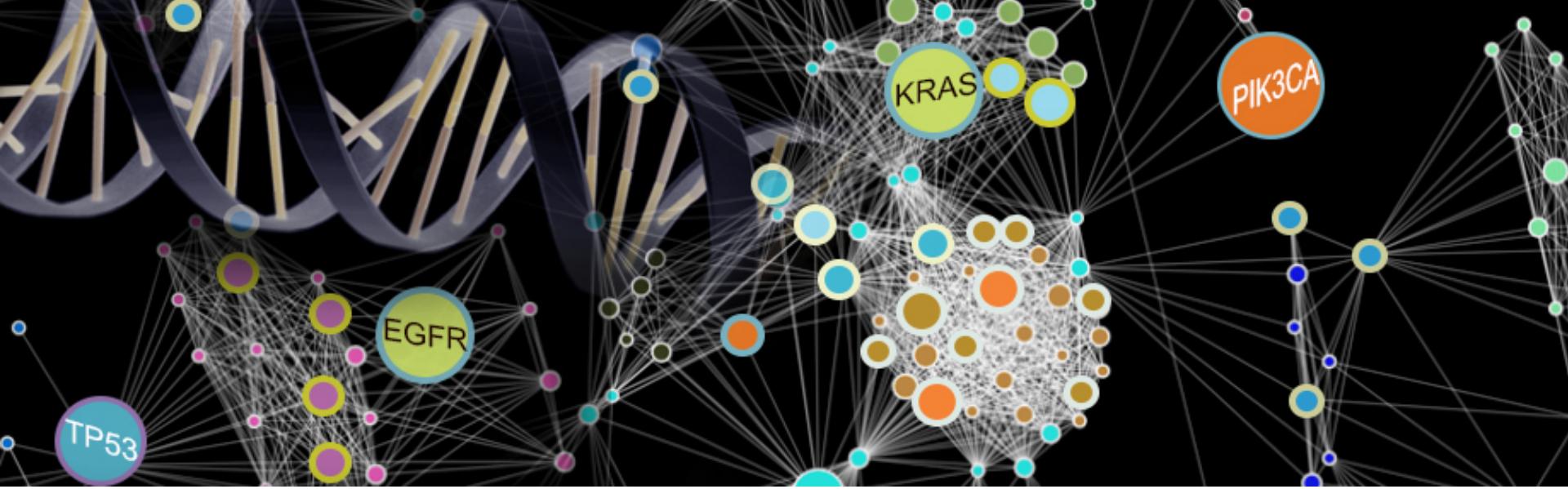
Tumor Type	Samples	Total Mutations	Missense	Silent	Nonsense	Nonstop	Splice	Frameshift	Muts/Samples
UCEC	542	492965	321267	100586	34237	403	10556	25916	909
SKCM	468	348935	212749	113524	16329	96	4525	1712	745
COAD	433	157925	95723	31742	8605	117	3049	18689	364
LUAD	567	174019	115051	38250	10162	165	4417	5974	306
STAD	439	133647	79156	28037	4767	108	2595	18984	304
LUSC	494	144624	95239	32962	7824	157	3581	4861	292
BLCA	412	107105	68343	26183	7785	138	1924	2732	259
READ	158	40260	27559	7134	3368	41	765	1393	254
CESC	305	64083	40664	15391	4003	93	1103	2829	210
GBM	396	63395	42546	13060	3449	60	1471	2809	160
HNSC	509	80761	51156	19274	4972	67	1686	3606	158
ESCA	184	25301	16519	5578	1155	27	505	1517	137
OV	443	60955	39682	11906	2962	62	1816	4527	137
PAAD	180	22487	14220	5522	1863	10	444	428	124
DLBC	48	5764	3742	1407	285	6	135	189	120
UCS	57	6800	4594	1238	505	9	137	317	119
LIHC	375	39321	25590	8804	1811	48	1136	1932	104
BRCA	1044	93317	58034	19406	6254	96	2360	7167	89
KIRP	288	23996	14740	5597	1047	30	521	2061	83
ACC	92	7200	4642	1568	422	7	170	391	78
SARC	255	17371	11369	4217	720	20	387	658	68
KIRC	339	21058	12161	4112	1290	22	577	2896	62
LGG	512	26084	16627	5886	1365	13	689	1504	50
LAML	145	5763	3363	1092	241	27	719	321	39

Mutation burden by tumor type



Lawrence et al 2013

Where do these mutations come from?

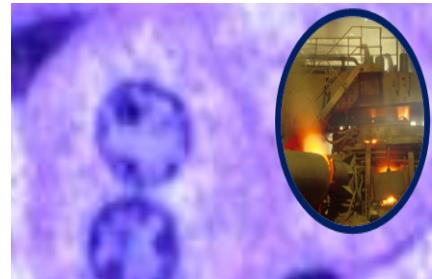


Mutation burden in tumors

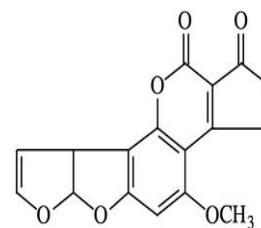
Why do tumor cells carry so many mutations?

Why is the burden different for different types of cancer?

Somatic mutations occur in all cells of the body throughout life



.....**A**T**C**GGG**A**AT**C**GG**A**CCC**G**AT**G**.....



Classification of base substitution mutations

C>T

C>A

C>G

T>A

T>C

T>G

Since DNA bases are paired
every possible substitution
can be summarized with 6
categories

By convention, these are
defined with respect to the
pyrimidine (C or T)

Classification of base substitution mutations

.....ATCGGGAAAT**C**GGACCCGATG.....
 ↓
.....ATCGGGAAAT**T**GGACCCGATG.....

Classification of base substitution mutations

.....ATCGGGAA**TCG**GACCCGATG.....
↓
.....ATCGGGAA**TTG**GACCCGATG.....

Classification of base substitution mutations

.....ATCGGGAA**TCG**GACCCGATG.....



.....ATCGGGAA**TTG**GACCCGATG.....

.....ATCGGGAA**ACG**GACCCGATG.....



.....ATCGGGAA**ATG**GACCCGATG.....

Classification of base substitution mutations

.....ATCGGGAA**TCG**GACCCGATG.....



.....ATCGGGAA**TTG**GACCCGATG.....

.....ATCGGGAA**ACG**GACCCGATG.....



.....ATCGGGAA**ATG**GACCCGATG.....

.....ATCGGGAA**ACC**GACCCGATG.....



.....ATCGGGAA**ATC**GACCCGATG.....

Classification of base substitution mutations

C>T

C>A

C>G

T>A

T>C

T>G

6 mutation classes

Classification of base substitution mutations

6 mutation classes

C>T

C>A

C>G

T>A

T>C

T>G

ACA>ATA

ACC>ATC

ACG>ATG

ACT>ATT

CCA>CTA

CCC>CTC

CCG>CTG

CCT>CTT

GCA>GTA

GCC>GTC

GCG>GTG

GCT>GTT

TCA>TTA

TCC>TTC

TCG>TTG

TCT>TTT

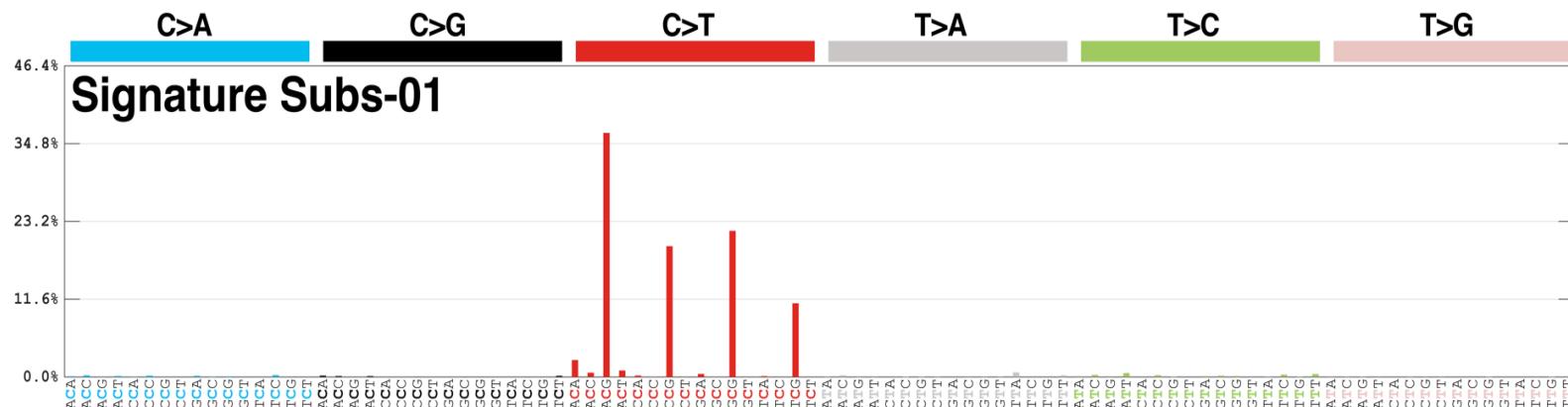
Classification of base substitution mutations

C>T	ACA>ATA ACC>ATC ACC>ATG ACT>ATT CCA>CTA CCC>CTC CCC>CTG CCG>CTT GCA>GTA GCC>GTC GCC>GTG GCF>GTT TCA>TTA TCC>TTC TCC>TTG TCT>TTT	ATA>AAA ATC>AAC ATG>AAG ATT>AAT CTA>CAA CTC>CAC CTG>CAG CTT>CAT GTA>GAA GTC>GAC GTG>GAG GTT>GAT TTA>TAA TTC>TAC TTG>TAG TTT>TAT
C>A	ACA>AAA ACC>AAC ACC>AAG ACT>AAT CCA>CAA CCC>CAC CCG>CAG CTC>CAT GCA>GAA GCC>GAC GCC>GAG GCI>GAT TCA>TAA TCC>TAC TCC>TAG TCT>TAT	ATA>ACA ATC>ACC ATG>ACG ATT>ACT CTA>CCA CTC>CCC CTG>CCG CTT>CCT GTA>GCA GTC>GCC GTG>GCG GTT>GCT TTA>TCA TTC>TCC TTG>TCG TTT>TCT
C>G	ACA>AGA ACC>AGC ACC>AGG ACI>AGT CCA>CGA CCC>CGC CCG>CGG CTC>CGT GCA>GGA GCC>GGC GCG>GGG GCI>GGT TCA>TGA TCC>TGC TCC>TGG TCT>TGT	ATA>AGA ATC>AGC ATG>AGG ATT>AGT CTA>CGA CTC>CGC CTG>CGG CTT>CGT GTA>GGA GTC>GGC GTG>GGG GTT>GGT TTA>TGA TTC>TGC TTG>TGG TTT>TGT
T>A		
T>C		
T>G		

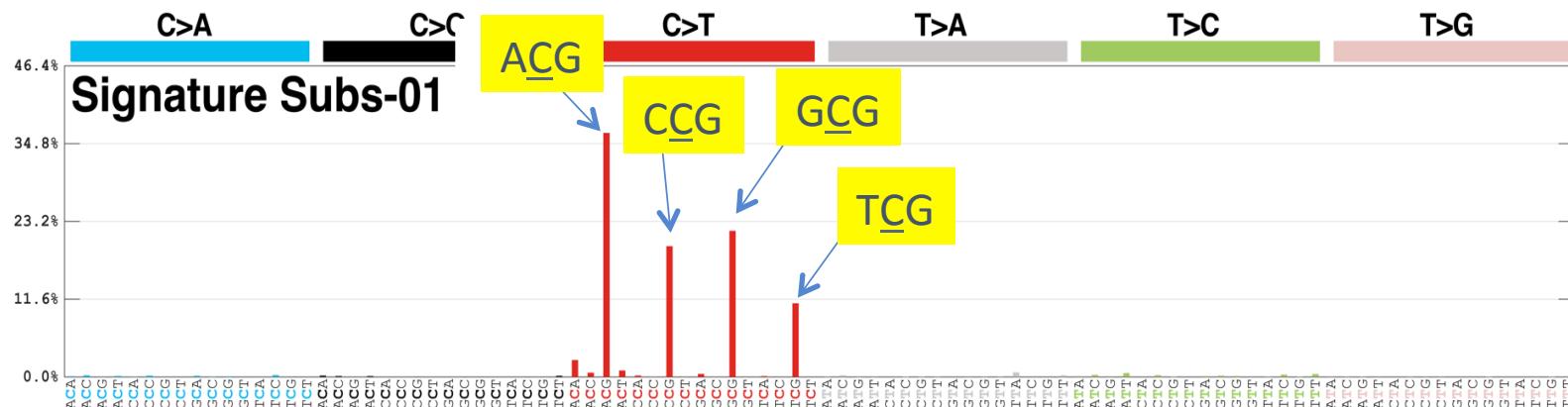
6 mutation classes

96 mutation classes

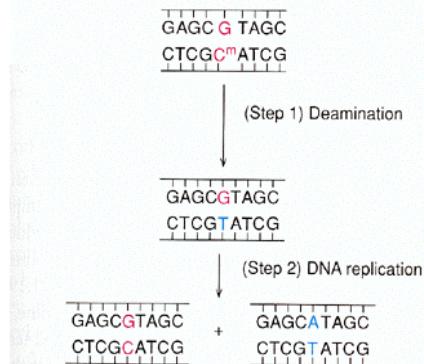
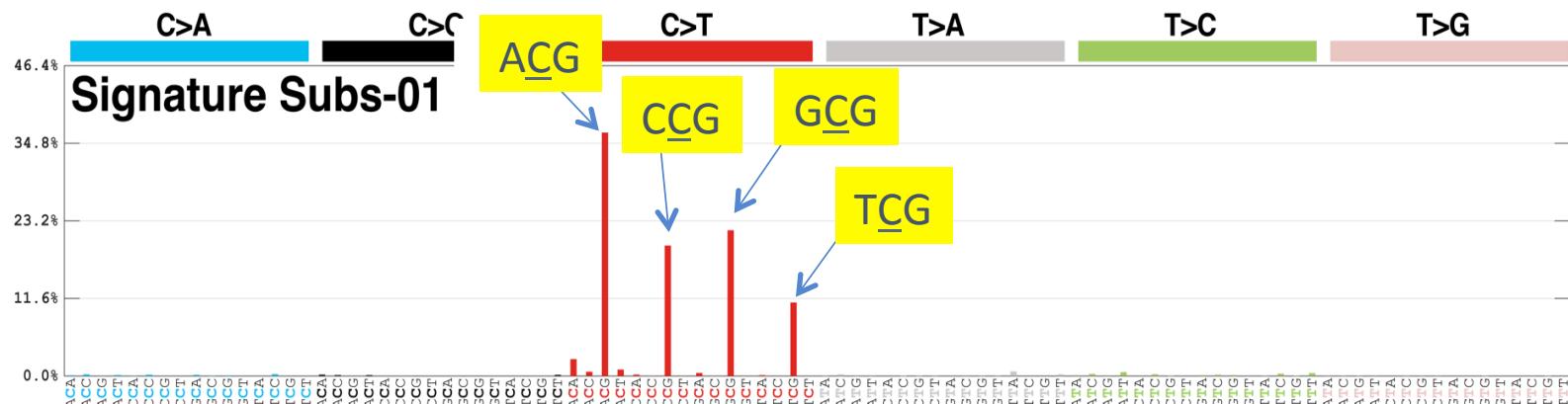
Mutational signatures in human cancer



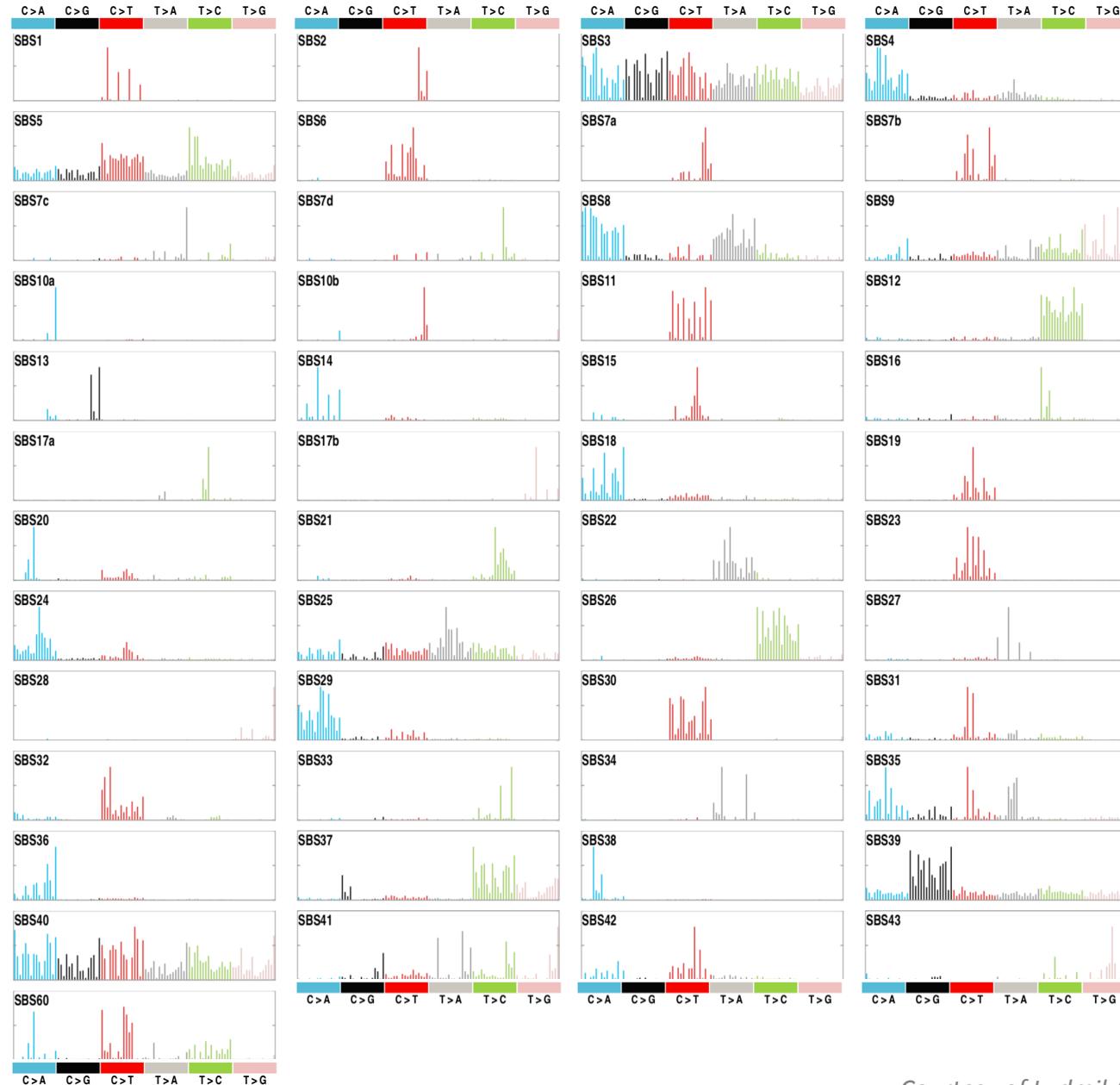
Mutational signatures in human cancer



Mutational signatures in human cancer

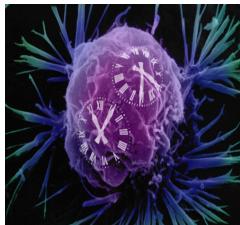


Aetiologies of substitution mutational signatures



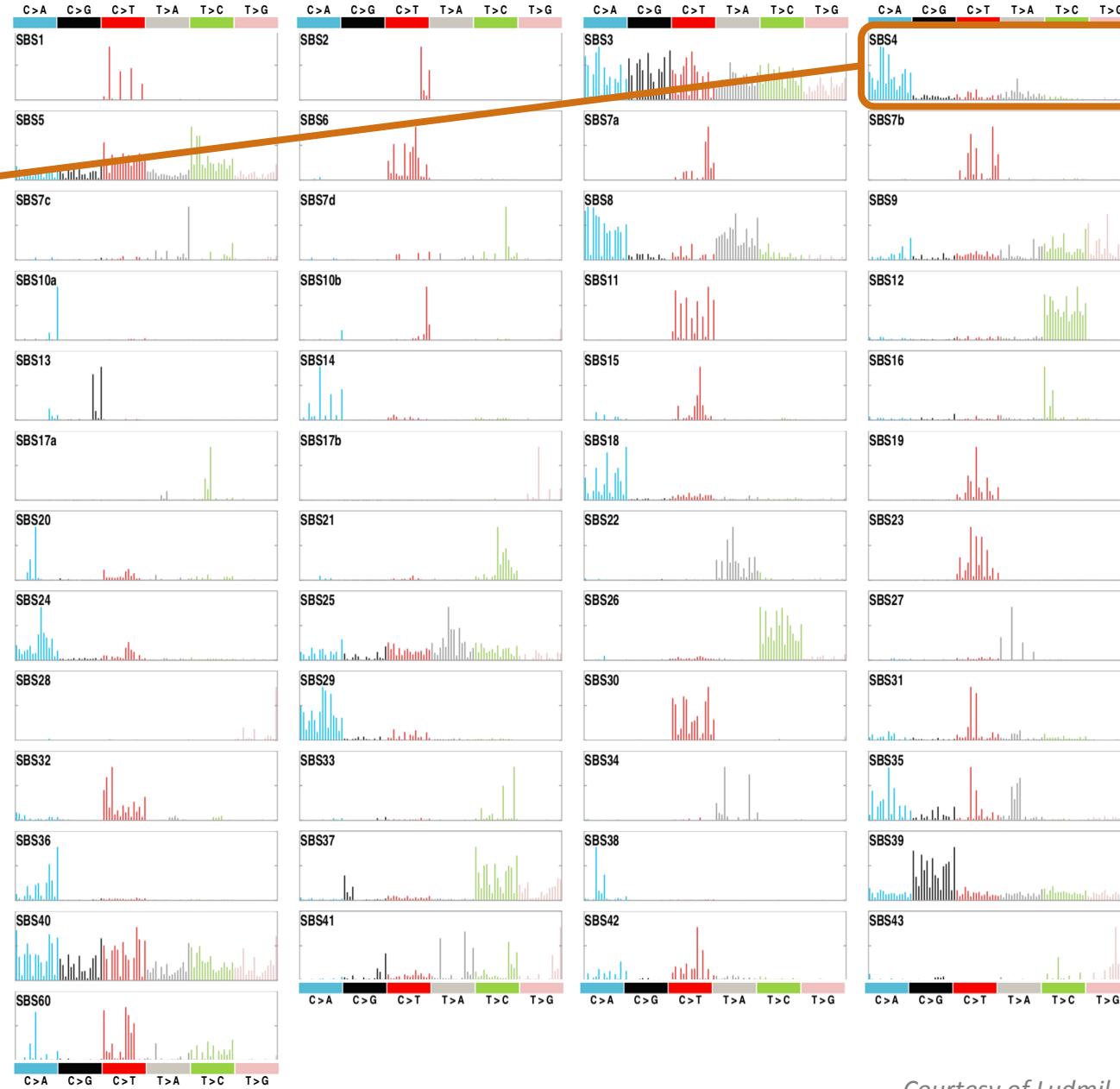
Aetiologies of substitution mutational signatures

Clock-like processes



Aetiologies of substitution mutational signatures

Tobacco smoking



Aetiologies of substitution mutational signatures

Tobacco chewing



Aetiologies of substitution mutational signatures

Ultraviolet
light



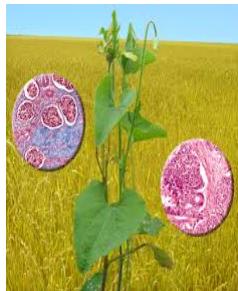
Aetiologies of substitution mutational signatures

Aflatoxin



Aetiologies of substitution mutational signatures

Aristolochic acid



Aetiologies of substitution mutational signatures

Temozolomide



Aetiologies of substitution mutational signatures

Platinum therapy



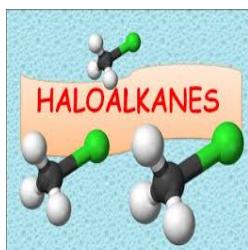
Aetiologies of substitution mutational signatures

Azathioprine



Courtesy of Ludmil Alexandrov

Aetiologies of **substitution** mutational signatures



Haloalkanes

Aetiologies of substitution mutational signatures

Defective
DNA
mismatch
repair



Aetiologies of substitution mutational signatures

Defective
BRCA1,
BRCA2,
homologous
recombination
repair



Aetiologies of substitution mutational signatures

Defective
base
excision
repair



Aetiologies of substitution mutational signatures

Defective polymerase epsilon activity



Aetiologies of substitution mutational signatures

Infidelity of polymerase eta activity



Aetiologies of substitution mutational signatures

APOBEC
cytosine
deamination



Aetiologies of substitution mutational signatures

Unknown
aetiologies
19/49



Origins of mutations observed in tumors

- DNA damaging environmental exposures
- Hereditary factors that impair the DNA damage response
- Replication error
 - Tissues with higher cell turnover rate will accumulate more mutations with age (Tomasetti and Vogelstein 2015)
 - The “bad luck” paper

⇒Mutations due to replication errors may account for more risk than hereditary and exposure related mutations
- Half of the mutations may pre-exist in cells prior to tumor initiation (Tomasetti et al 2012)

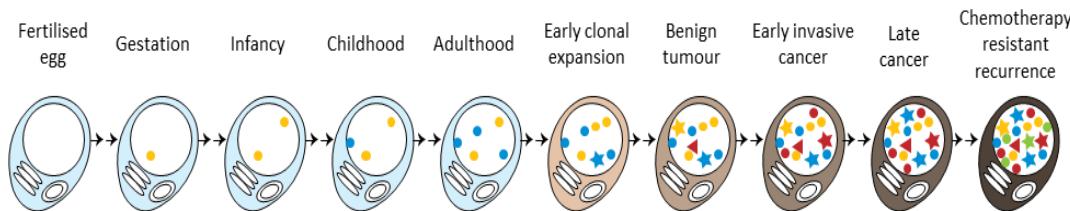
How do these many mutation sources contribute to the mutation burden of a tumor?

From fertilised egg to cancer cell

Chemotherapy
resistant
recurrence



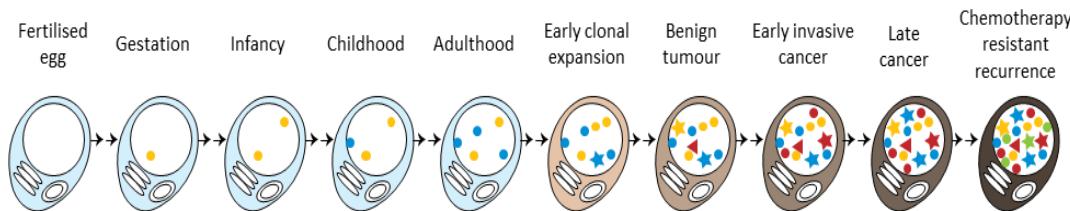
From fertilised egg to cancer cell



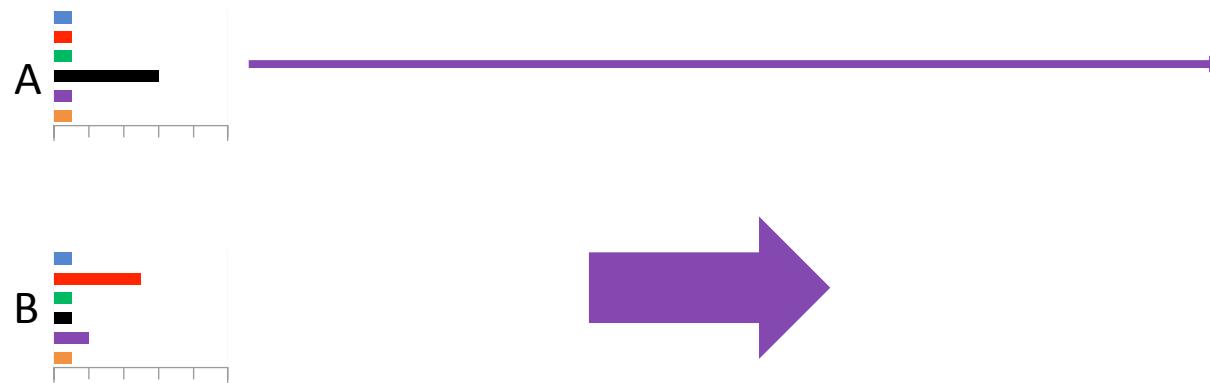
Mutational
processes



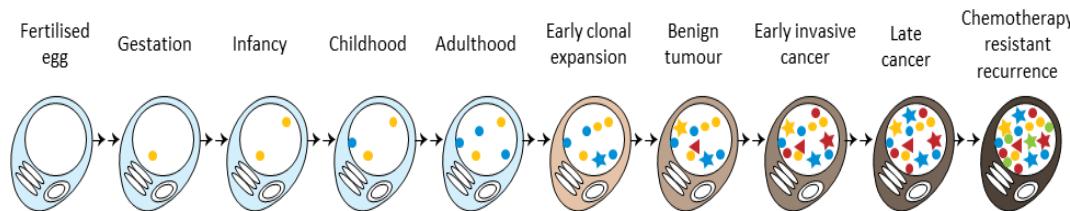
From fertilised egg to cancer cell



Mutational
processes



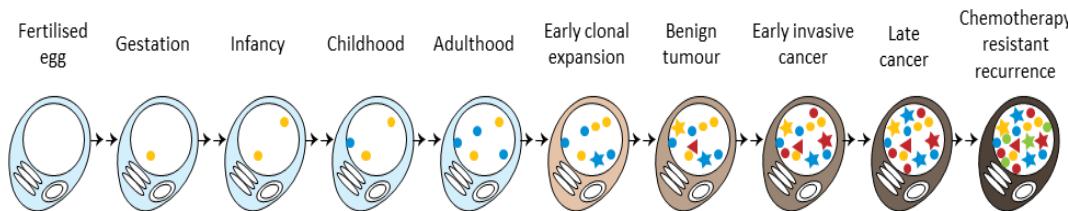
From fertilised egg to cancer cell



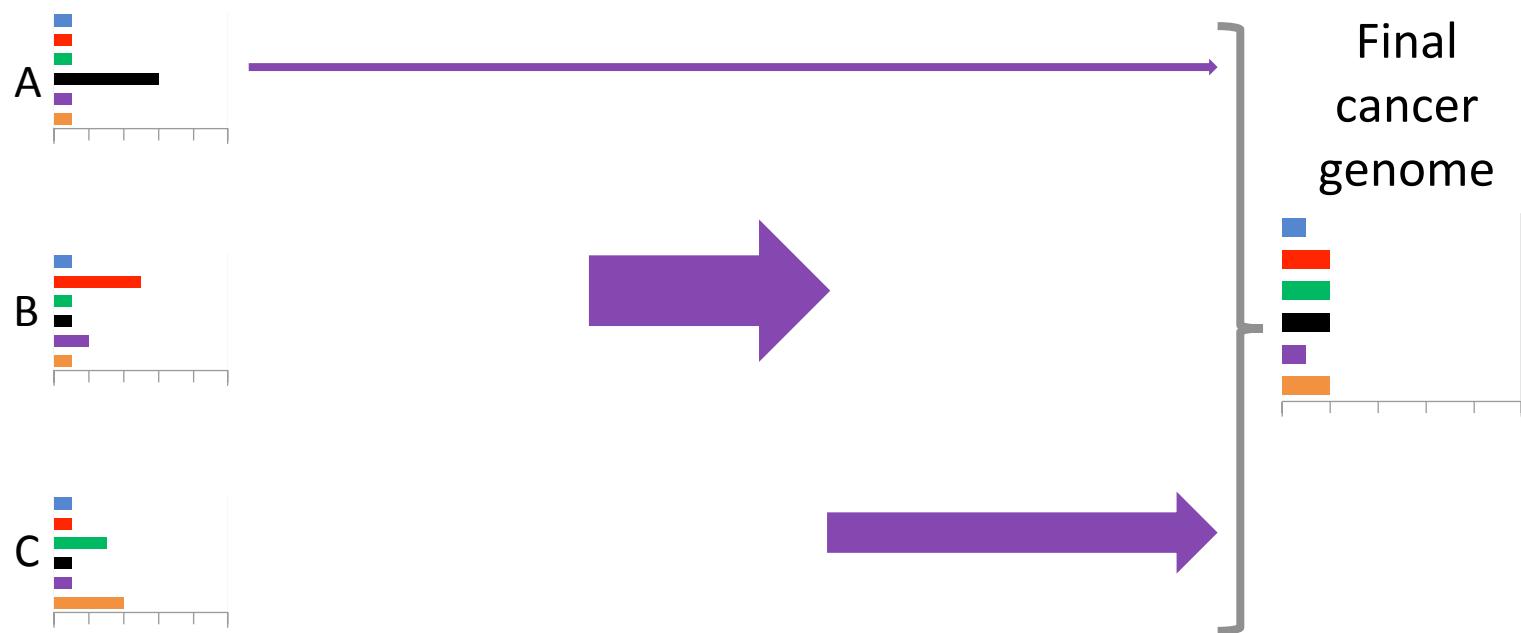
Mutational
processes



From fertilised egg to cancer cell

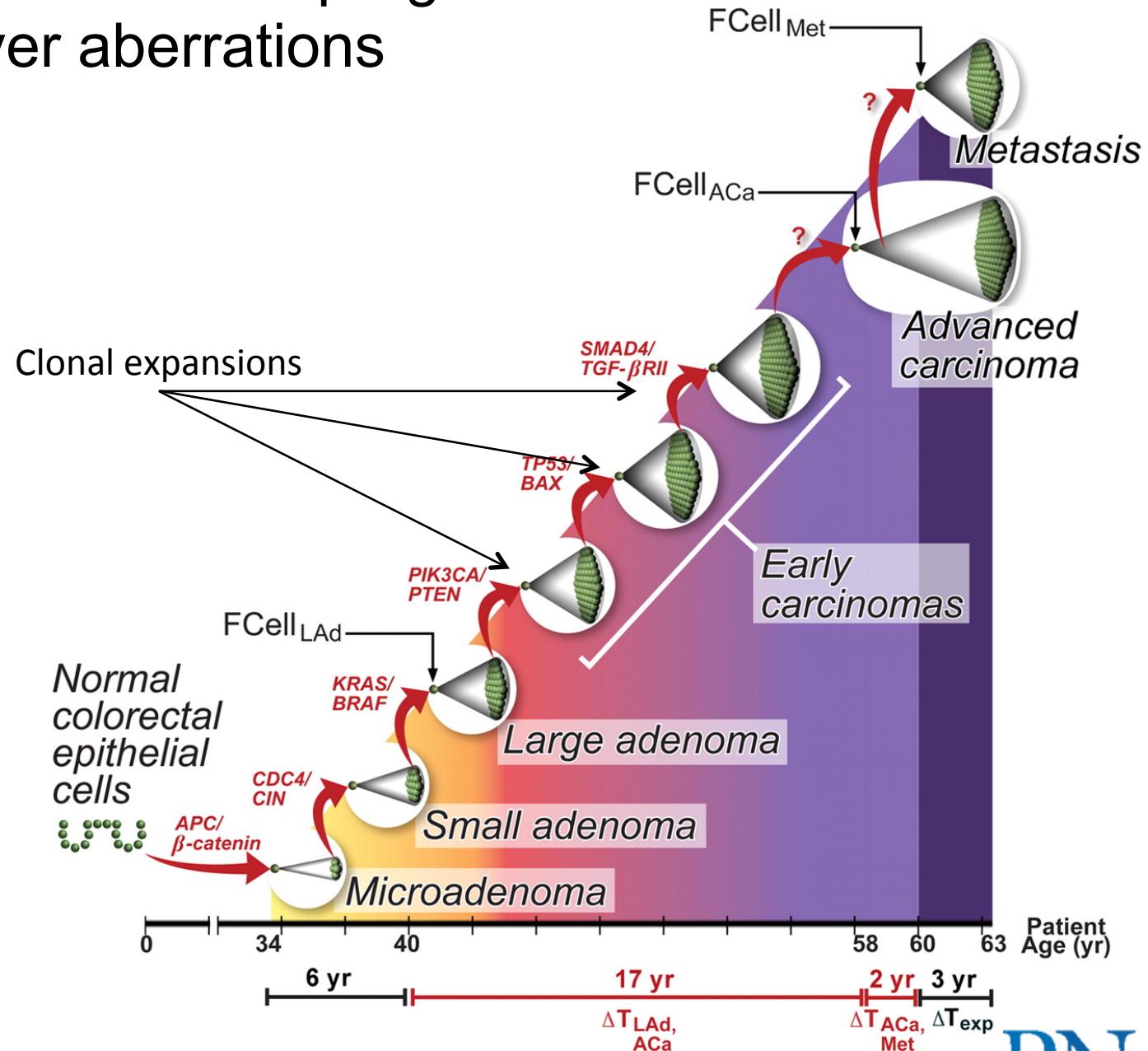


Mutational
processes



How does mutation burden relate to cancer risk and tumor development?

Cancer is a disease of progressive genetic driver aberrations

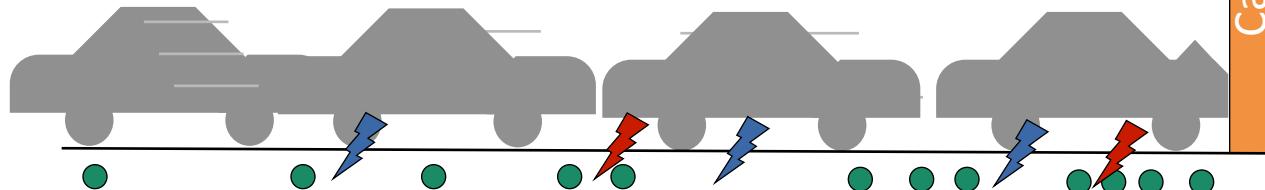


Driver versus Passenger mutations

Driver Mutations



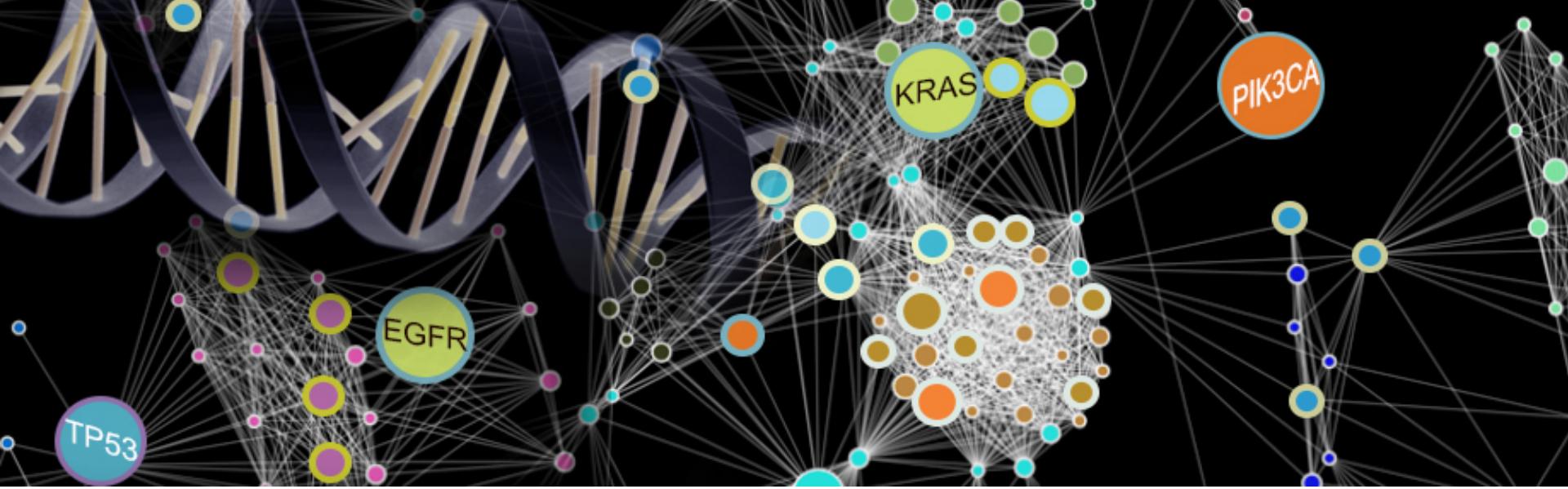
Mutations in oncogenes and tumor suppressors



Passenger Mutations

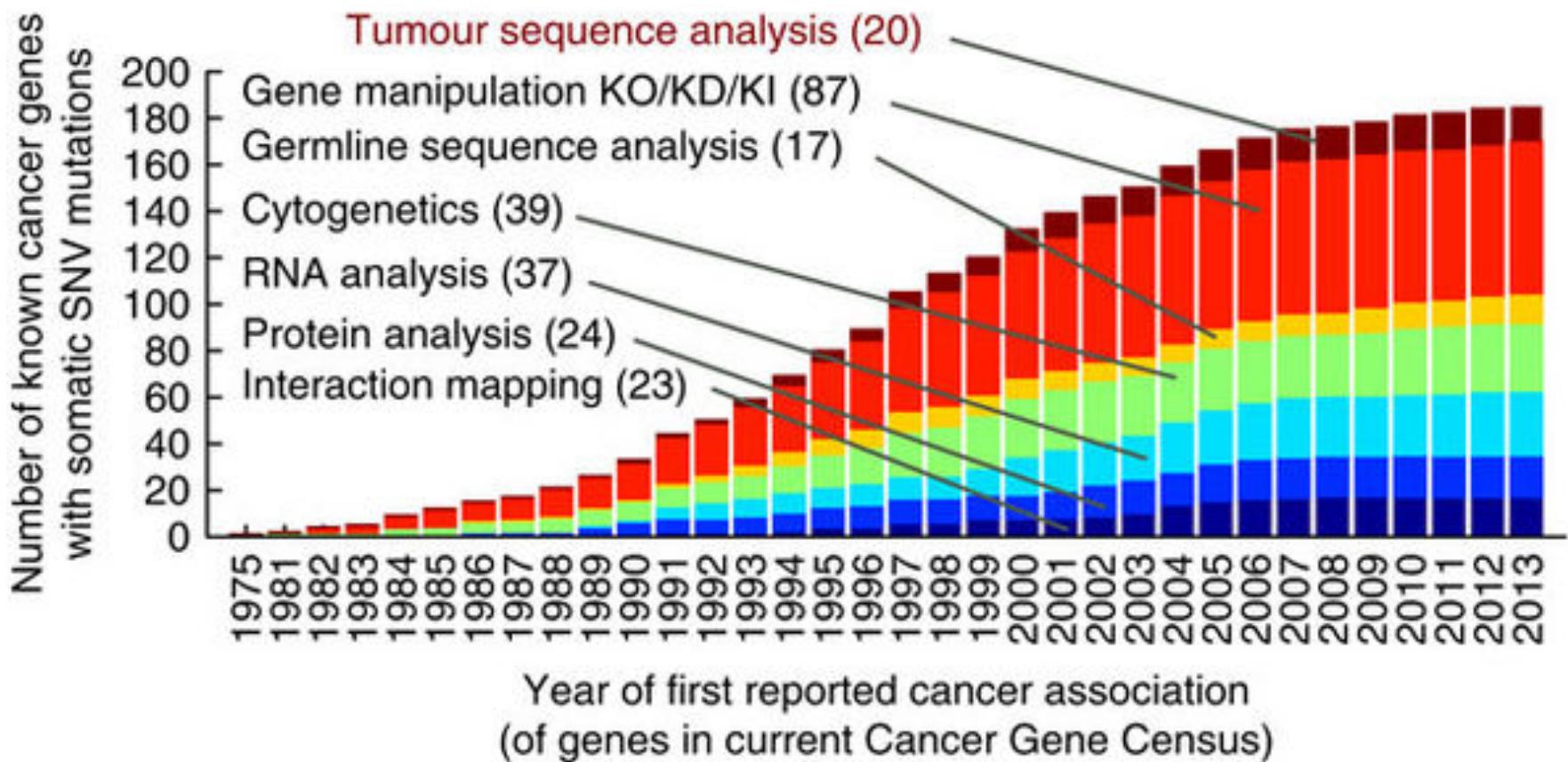


- Driver mutations provide a growth / survival advantage
- Passenger mutations can be preexist in a cell that becomes tumorigenic
- Passengers continue to accumulate after cancer initiation



Identifying driver mutations in cancer

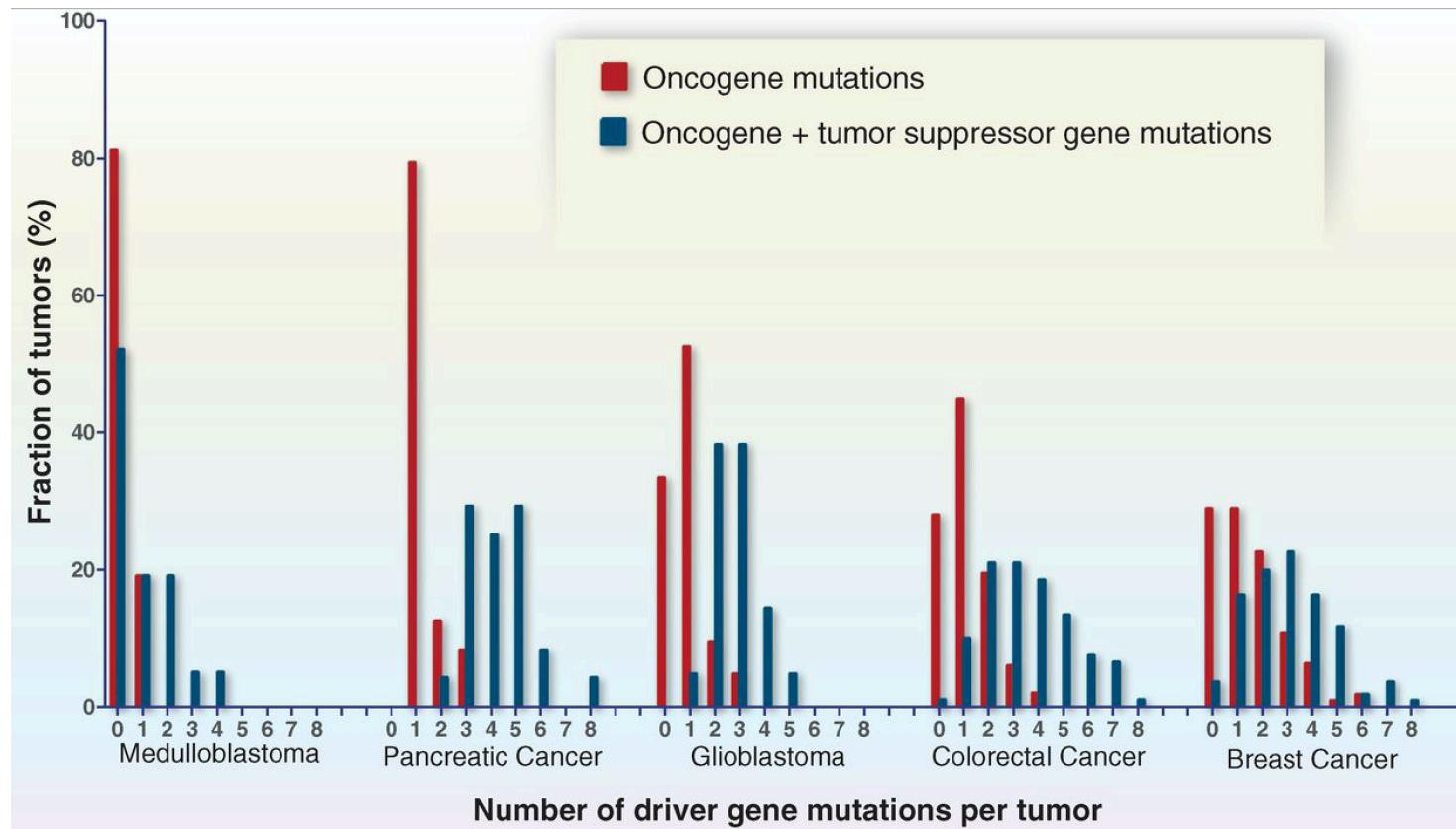
Most known driver genes were found before NGS



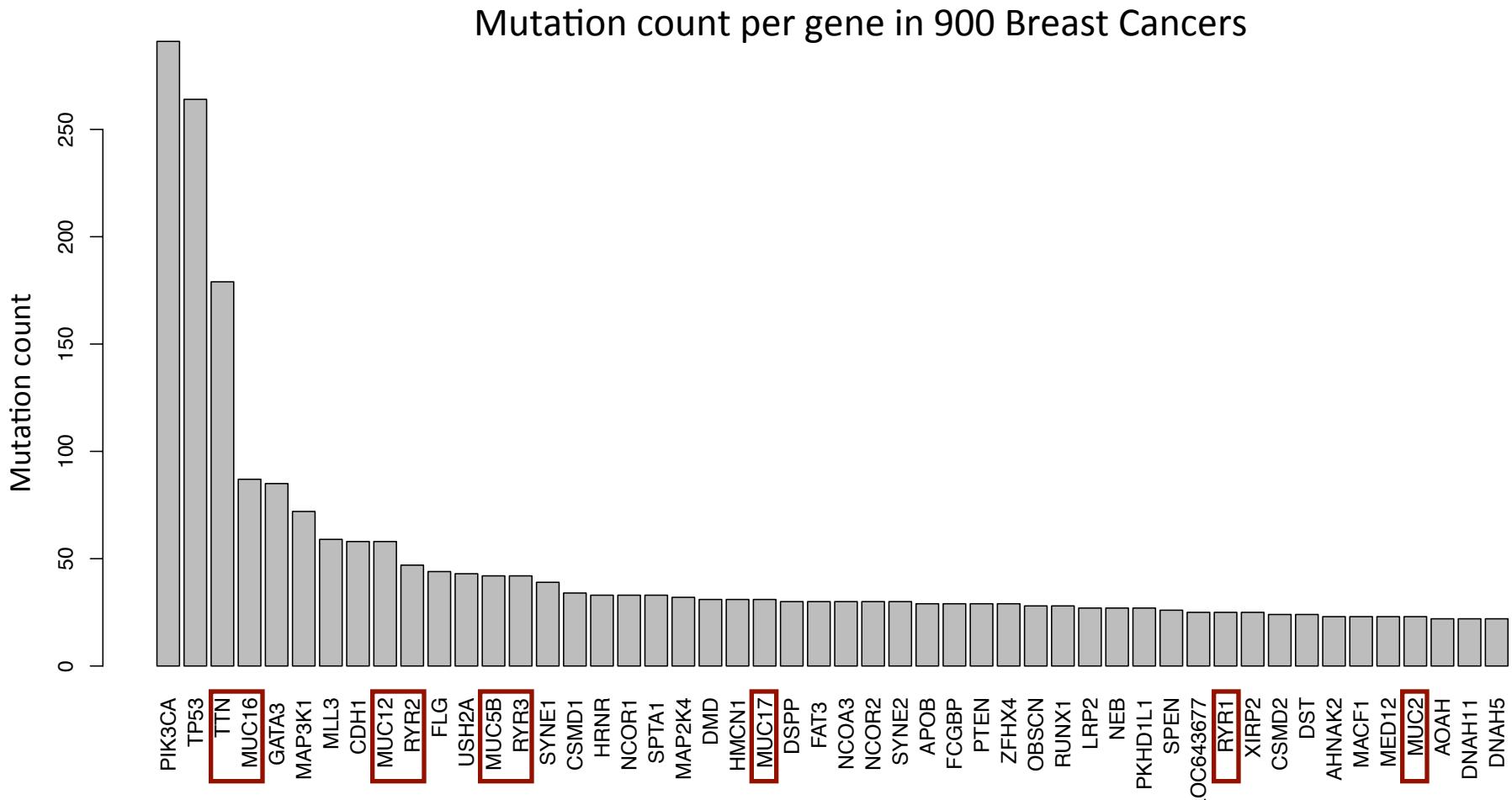
Most mutations in tumors are “Passengers”

Early epidemiological studies suggested ~2-8 drivers per tumor

NGS studies



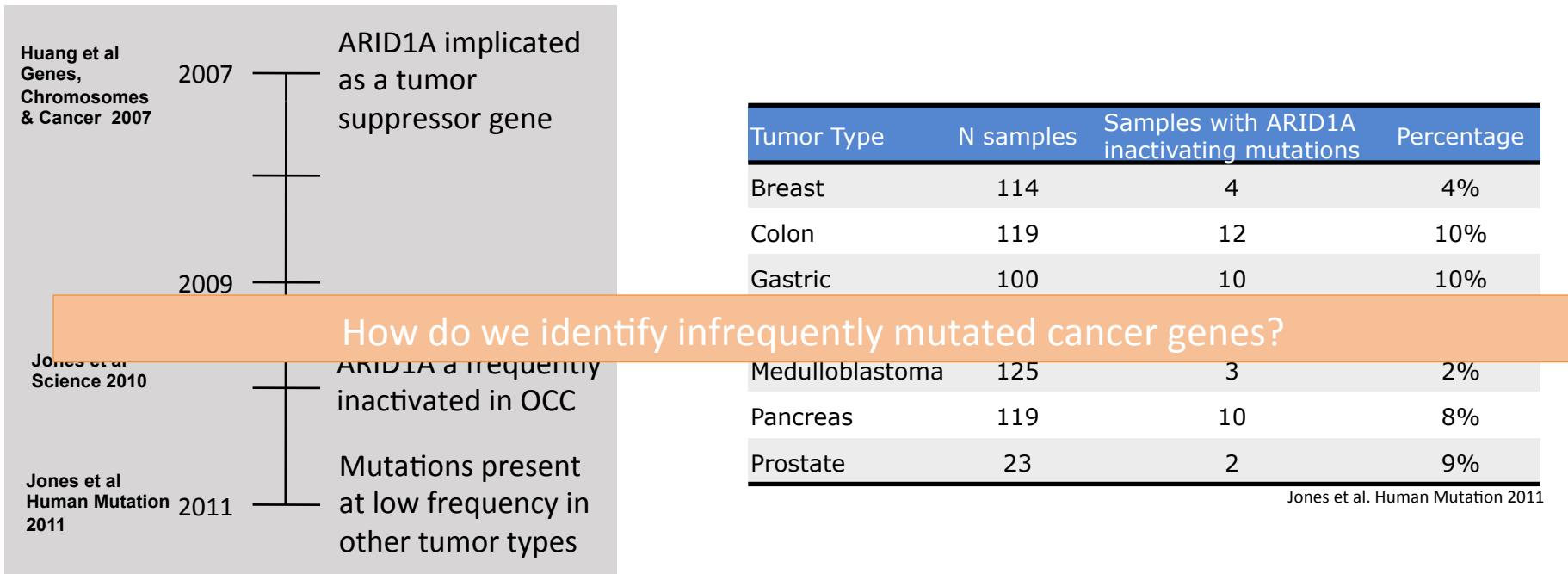
Frequently mutated genes



Frequency not a perfect indicator of a cancer gene

Evidence suggests *infrequently* mutated genes may harbor drivers

- ARID1A



If genes mutated at low frequency include known cancer genes, they may also include unknown cancer genes.

Existing Tools For Variant Prioritization

Name	Type	Information	URL	Refs
MAPP	Constraint-based predictor	Evolutionary and biochemical	http://mendel.stanford.edu/SidowLab/downloads/MAPP/index.html	27
SIFT	Constraint-based predictor	Evolutionary and biochemical (indirect)	http://sift.bii.a-star.edu.sg/	39
PANTHER	Constraint-based predictor	Evolutionary and biochemical (indirect)	http://www.pantherdb.org/	41
MutationTaster*	Trained classifier	Evolutionary, biochemical and structural	http://www.mutationtaster.org/	40
nsSNP Analyzer	Trained classifier	Evolutionary, biochemical and structural	http://snpanalyzer.uthsc.edu/	44
PMUT	Trained classifier	Evolutionary, biochemical and structural	http://mmbr2.pcb.ub.es:8080/PMut/	38
polyPhen	Trained classifier	Evolutionary, biochemical and structural	http://genetics.bwh.harvard.edu/pph2/	35
SAPRED	Trained classifier	Evolutionary, biochemical and structural	http://sapred.cbi.pku.edu.cn/	42
SNAP	Trained classifier	Evolutionary, biochemical and structural	http://www.rostlab.org/services/SNAP/	36
SNPs3D	Trained classifier	Evolutionary, biochemical and structural	http://www.snps3d.org/	51
PhD-SNP	Trained classifier	Evolutionary and biochemical (indirect)	http://gpcr2.biocomp.unibo.it/~emilio/PhD-SNP/PhD-SNP_Help.html	37

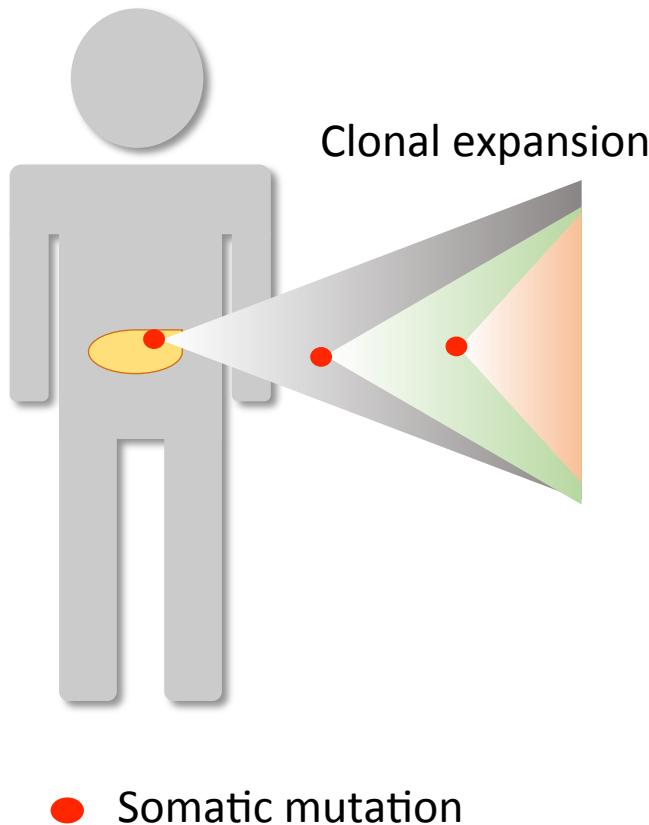
*Also makes predictions for synonymous and non-coding variant effects: for example, splicing. MAPP, Multivariate Analysis of Protein Polymorphism; polyPhen, polymorphism phenotyping.

And many, many more ...

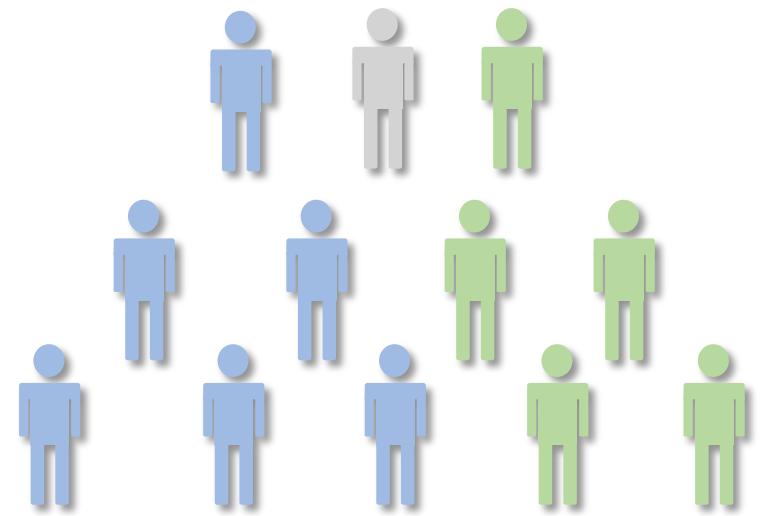
Cooper, G. M. & Shendure, J. Needles in stacks of needles: finding disease-causal variants in a wealth of genomic data. *Nat Rev Genet* **12**, 628–640 (2011).

Why do we need a new tool for cancer?

Accumulation of Mutations

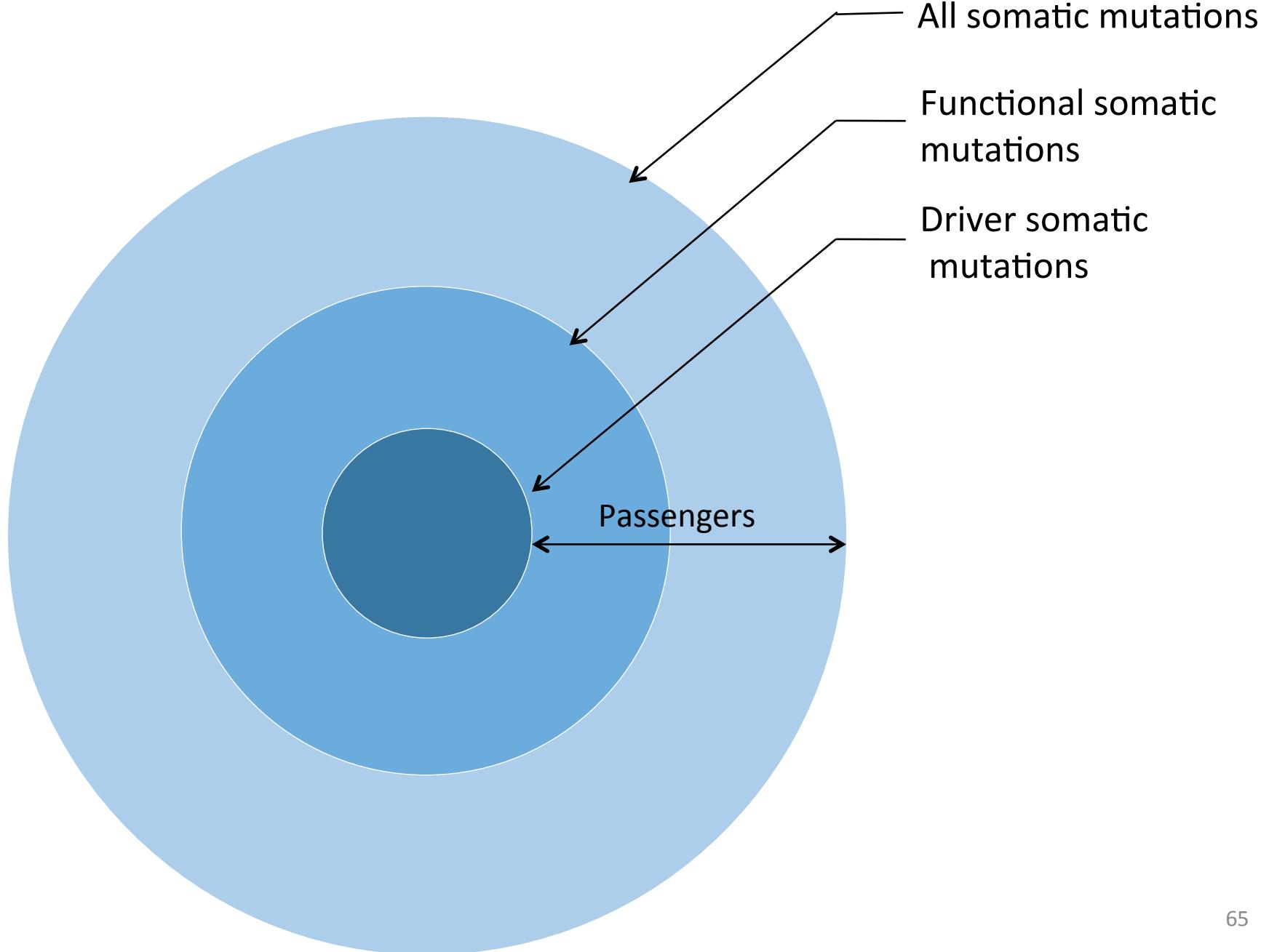


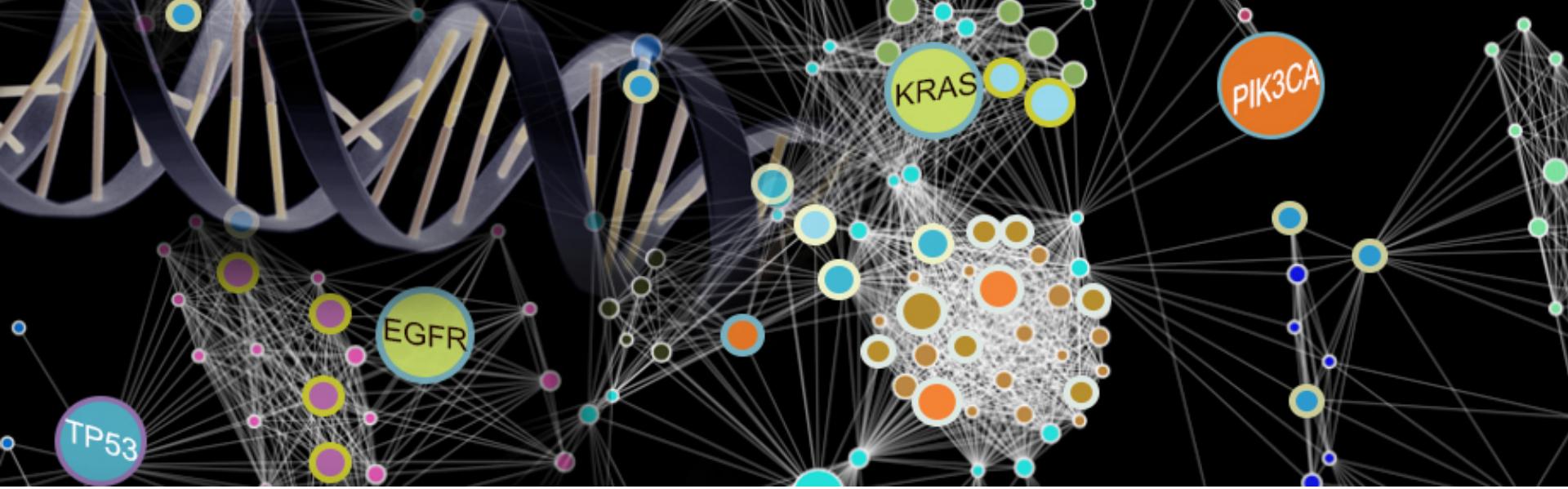
Germline evolution



- Rare inherited disease mutation
- Common inherited disease mutation
- Common neutral polymorphism

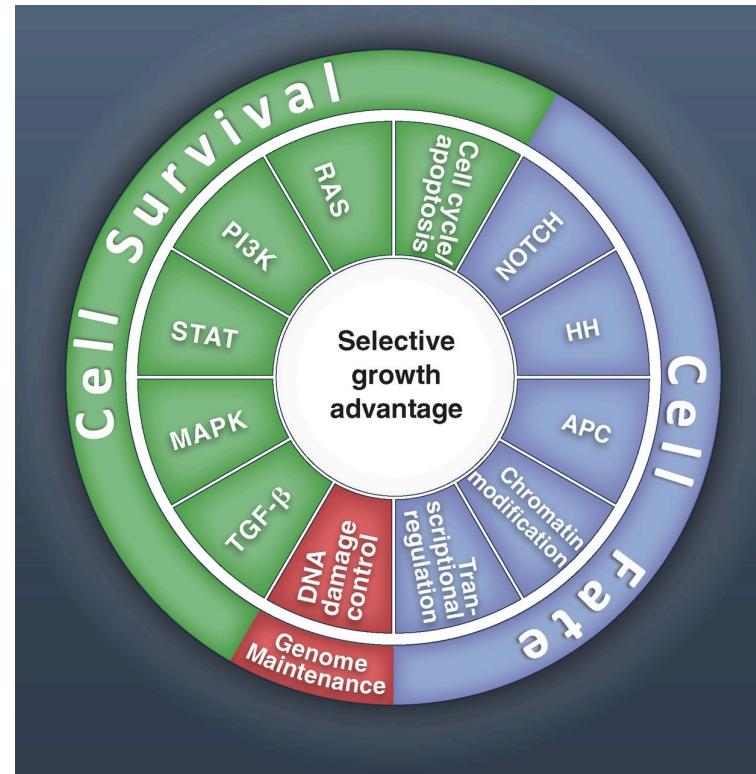
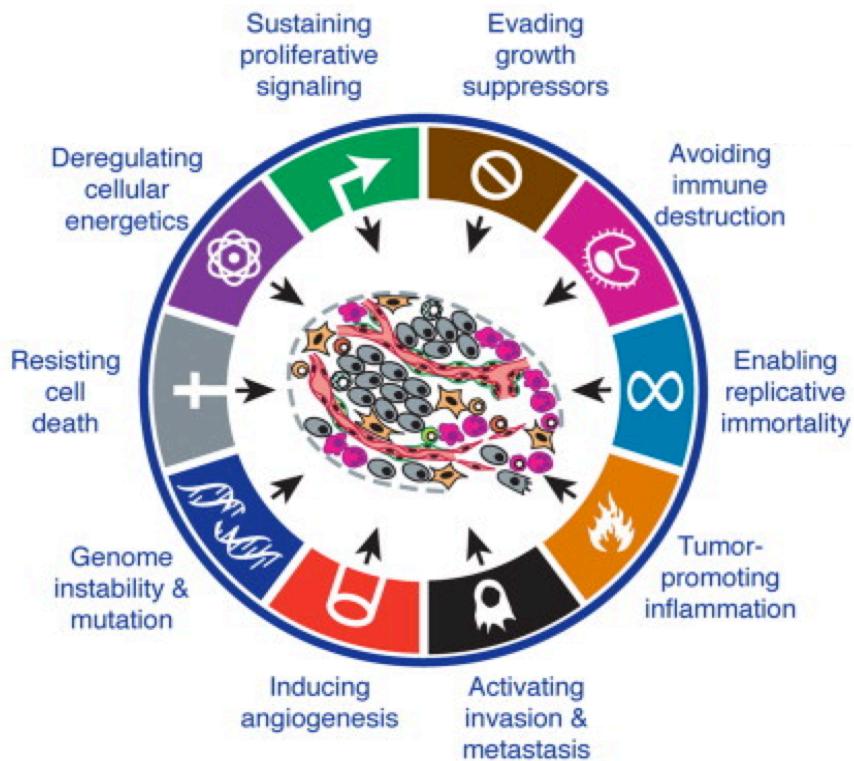
Cancer evolution is different from germline evolution





Tumor heterogeneity

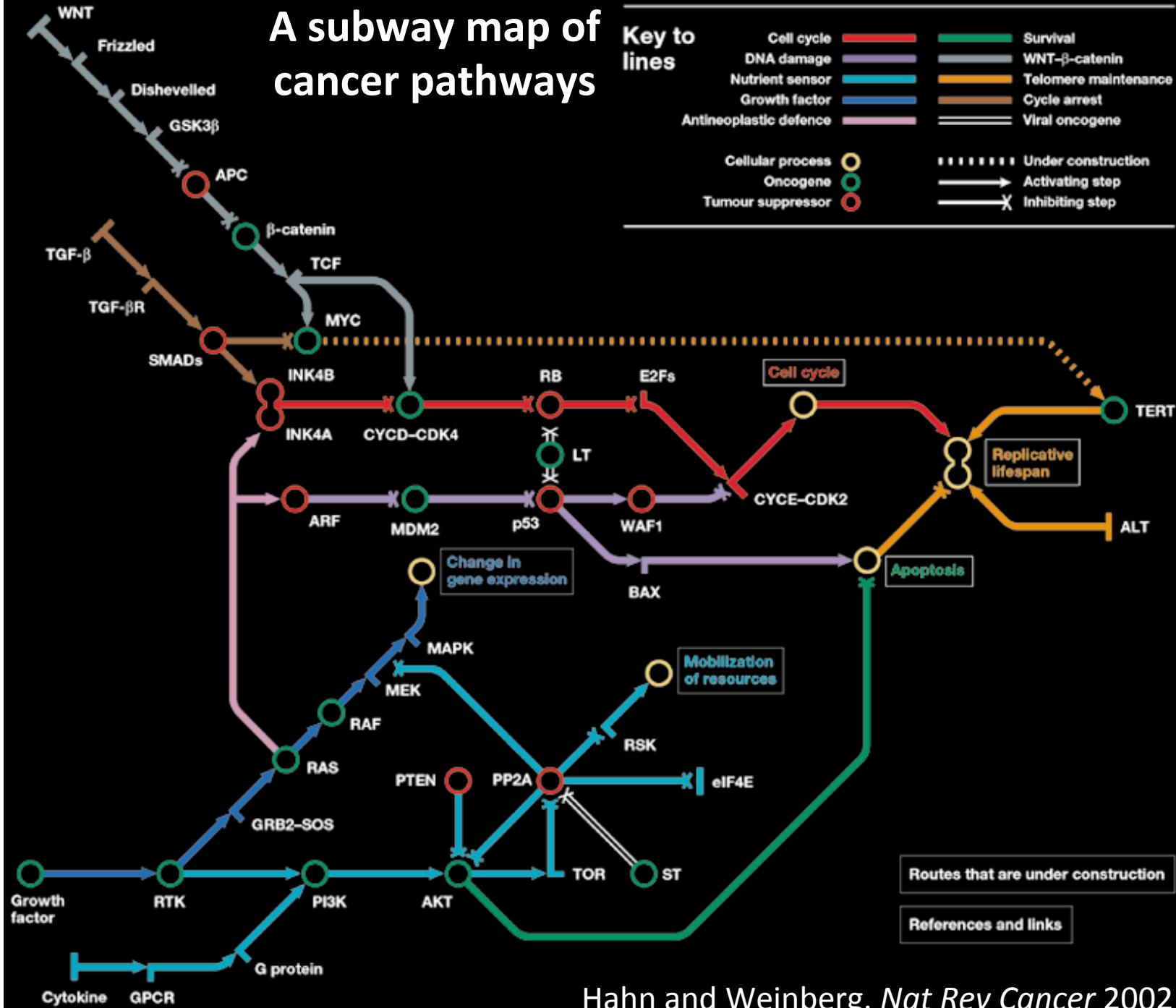
Cancer: Different mutations, same phenotypic characteristics



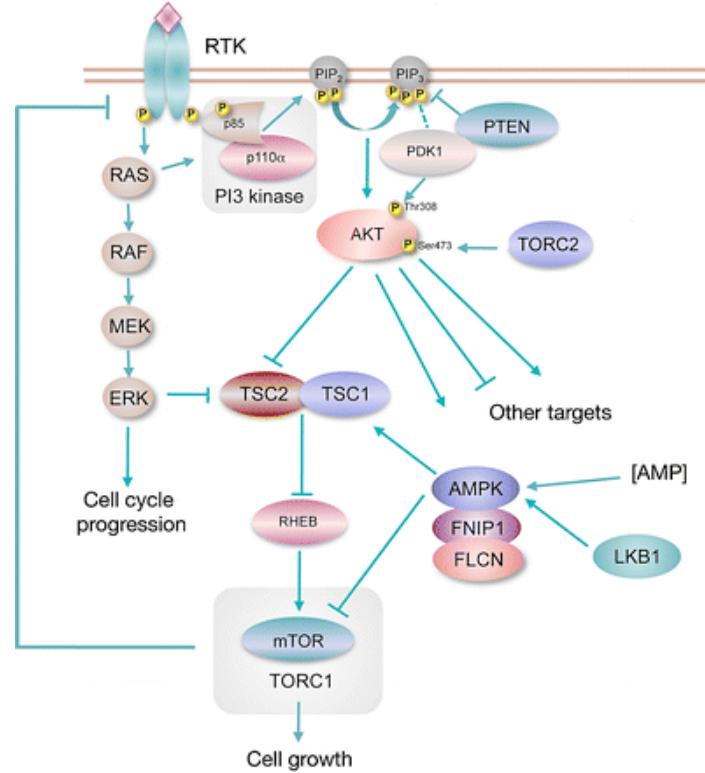
Hanahan, Weinberg **Hallmarks of Cancer: The Next Generation** Cell
Volume 144, Issue 5, 2011, 646 – 674

B Vogelstein et al. **Cancer Genome Landscapes** Science
2013, 339, 1546-1558

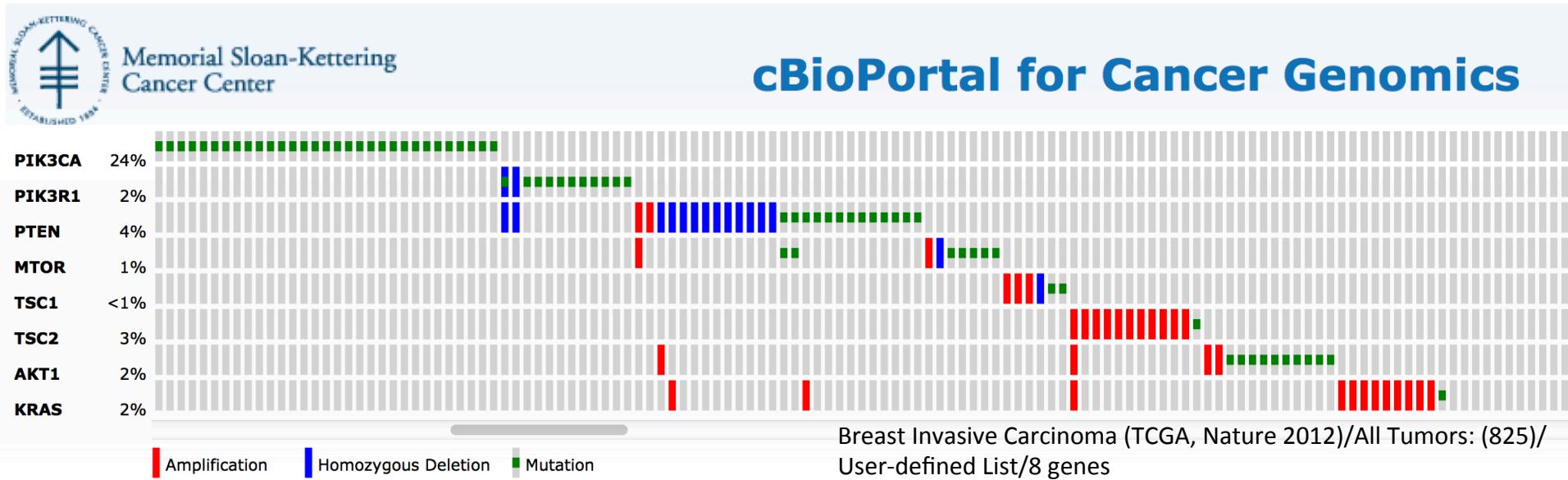
A subway map of cancer pathways



Mutual exclusivity within pathways



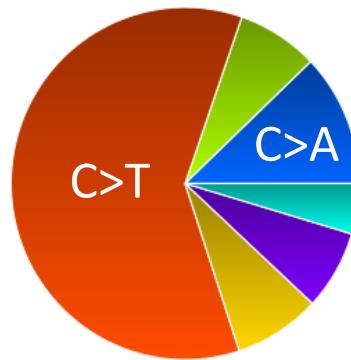
Knowles *et al* Phosphatidylinositol 3-kinase (PI3K) pathway activation in bladder cancer Cancer and Metastasis Reviews 2009



Tissue-specific mutation signatures can cause tissue-specific heterogeneity in drivers

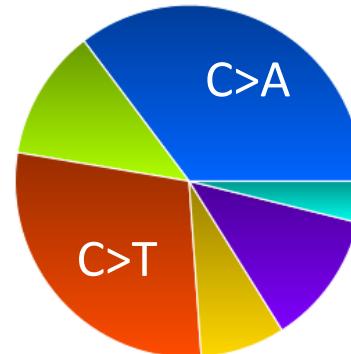
TP53: Most frequently mutated driver gene in cancer

686 skin cancers



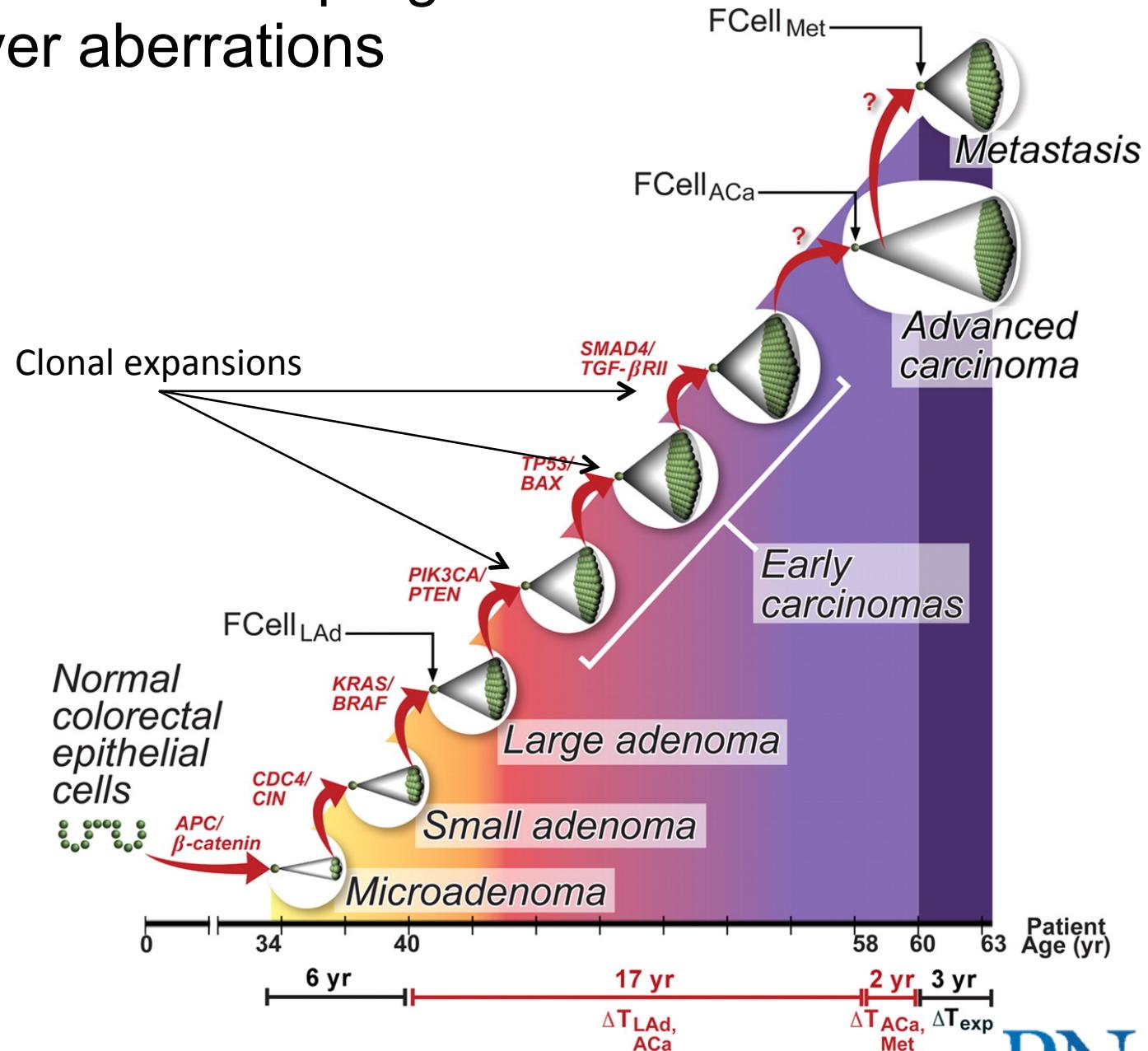
Ultraviolet light causes C>T mutations

1647 lung cancers



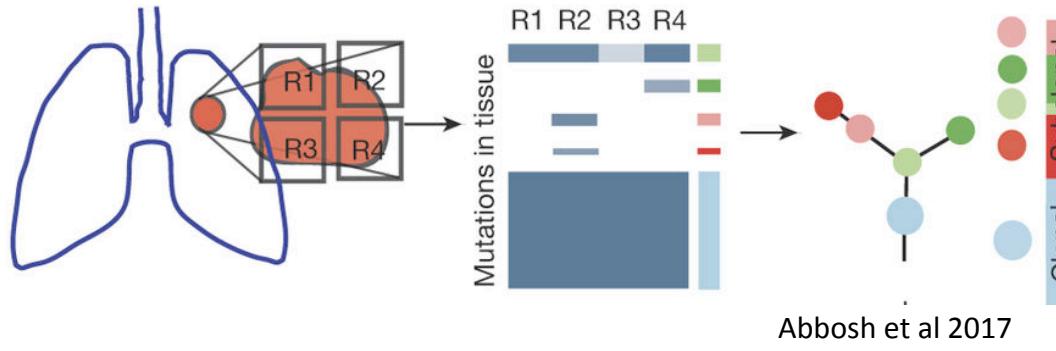
Tobacco carcinogens cause C>A mutations

Cancer is a disease of progressive genetic driver aberrations



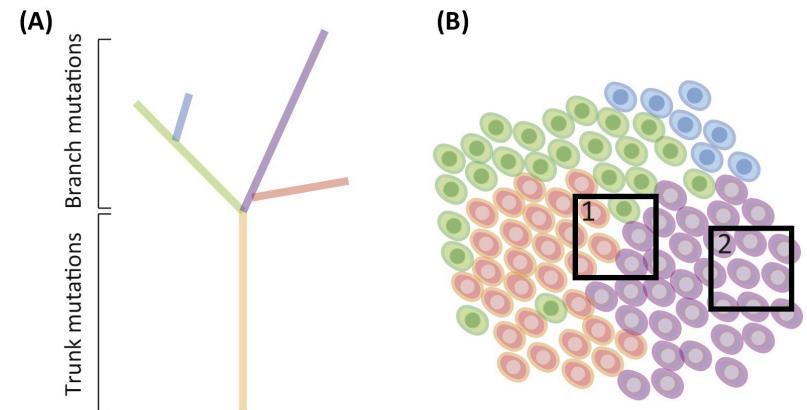
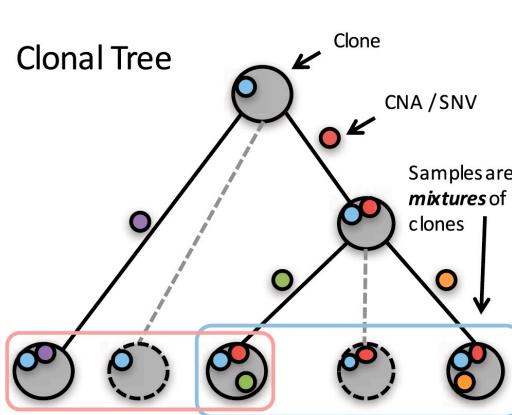
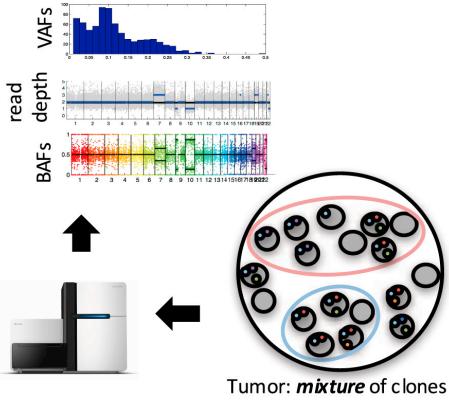
Clonal expansion creates spatial heterogeneity of mutations

- Sequencing different locations reveals different mutations



Abbosh et al 2017

- Sequencing coverage of mutations can be used to create hypothetical models of clonal composition but is challenged by sequencing bias



Mutation burden in tumors

More than just a count of mutations

Contains the history of DNA damage experienced by the tumor cell lineage

Not all mutations co-exist in the same tumor cells

Subclonal heterogeneity a major barrier for successful treatment

Need to eliminate all tumor cells but therapies may only target a subset

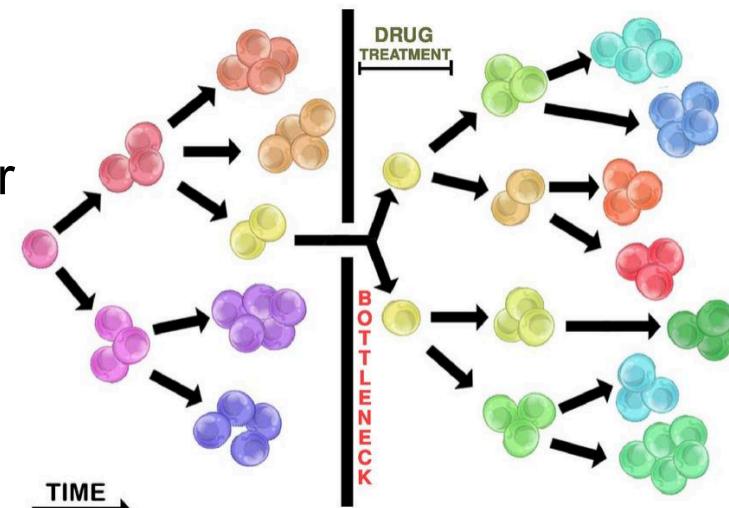
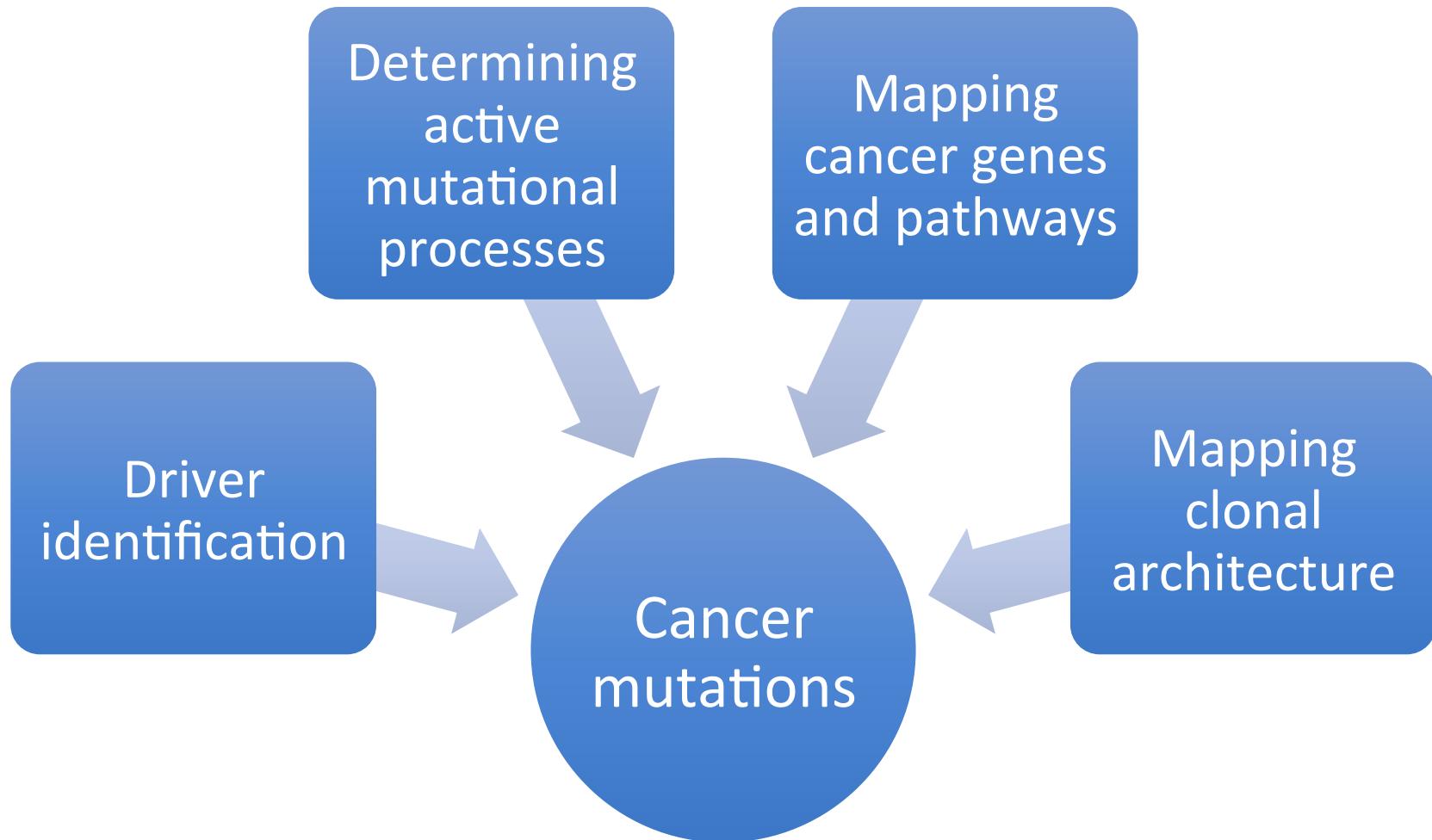
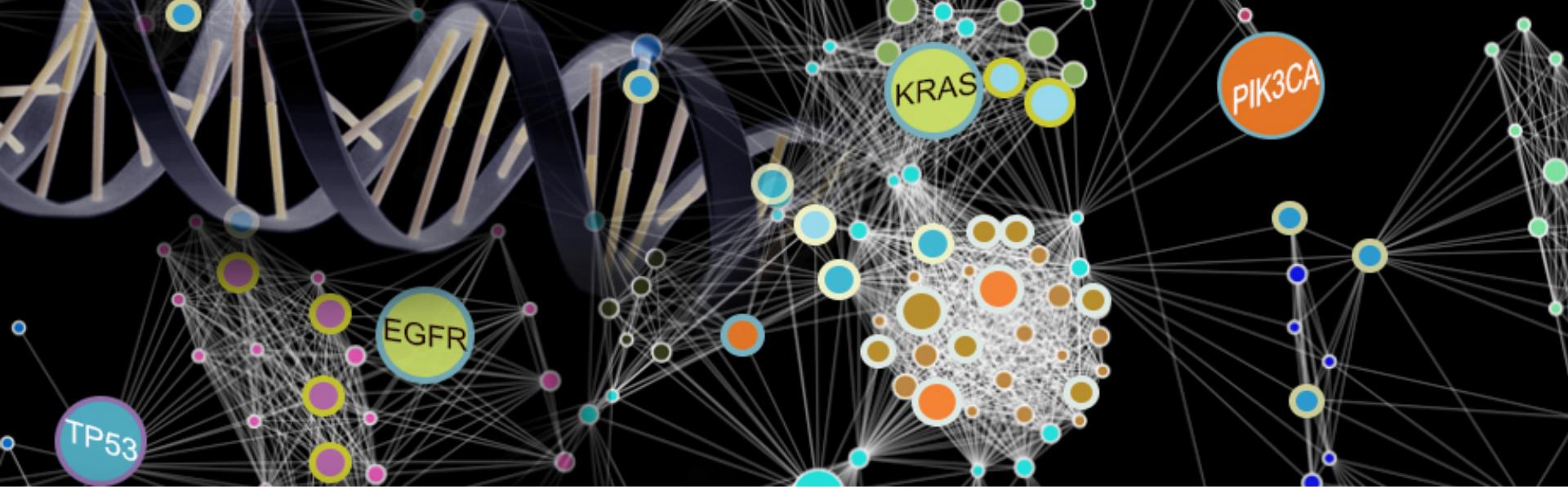


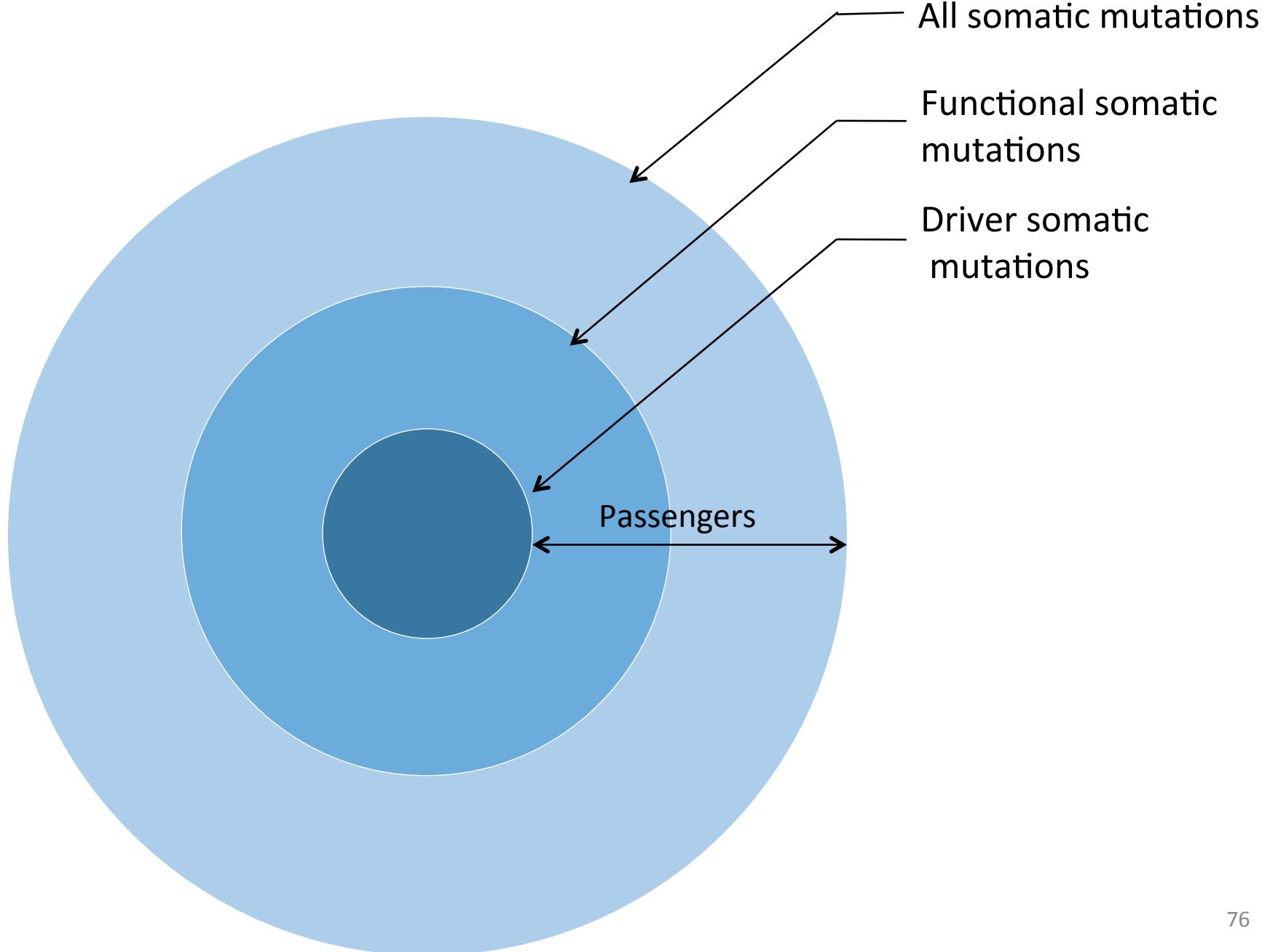
Image source: Wikipedia page on tumor heterogeneity

Major bioinformatic challenges for studying cancer mutations





Practical Example: Machine learning for driver mutation identification



Machine Learning: Supervised Prioritization

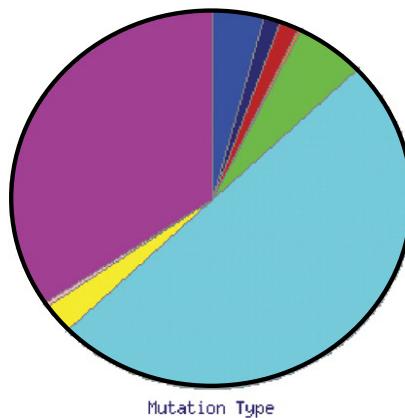
- **Training Set**
 - Positive examples (drivers)
 - Negative examples (passengers)
- Feature set
 - Informative characteristics of amino acid substitutions

Positive class: Driver Mutations



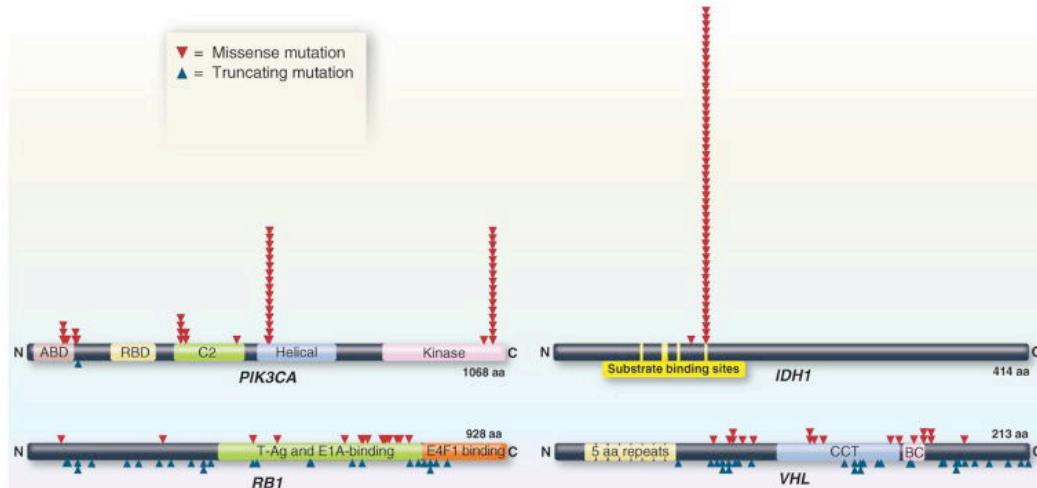
Curated data type	Curated data count
Experiments	2,760,220
Tumours	541,928
Mutations	136,326
References	10,383
Genes	18,490
Fusions	4946
Structural variants	2307
Whole cancer genomes	29

Forbes S A et al. Nucl. Acids Res. 2011;39:D945-D950



Details For Mutation Type Chart				
	Mutation Type	Number	Percentage	Mutation Data
1	Deletion Frameshift	855	4.2 %	More Details
2	Deletion Inframe	275	1.3 %	More Details
3	Insertion Frameshift	317	1.5 %	More Details
4	Insertion Inframe	50	0.2 %	More Details
5	Substitution Nonsense	1139	5.5 %	More Details
6	Substitution Missense	10292	50.0 %	More Details
7	Substitution Synonymous	518	2.5 %	More Details
8	Complex	65	0.3 %	More Details
9	Other	7056	34.3 %	More Details
10	Total	20567	100%	More Details

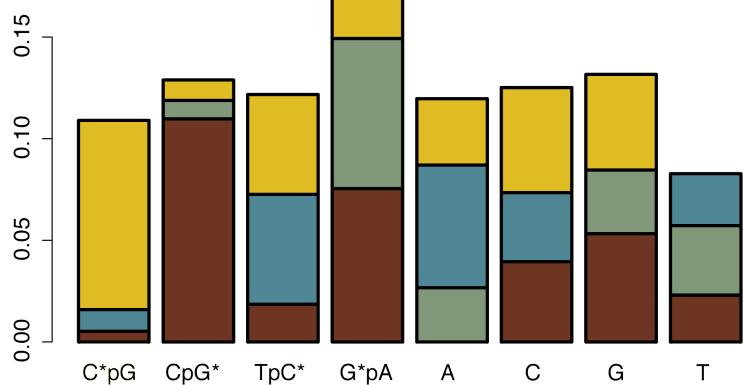
Selection Criteria:



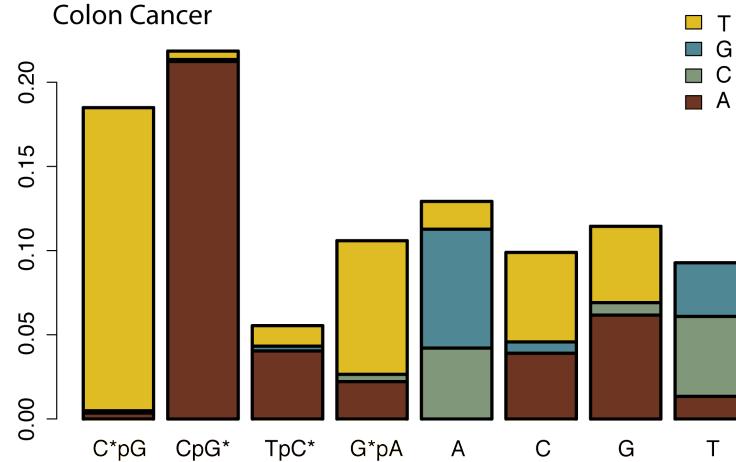
What do passenger mutations look like?

Assumption: The majority of mutations detected in tumor sequencing are passengers

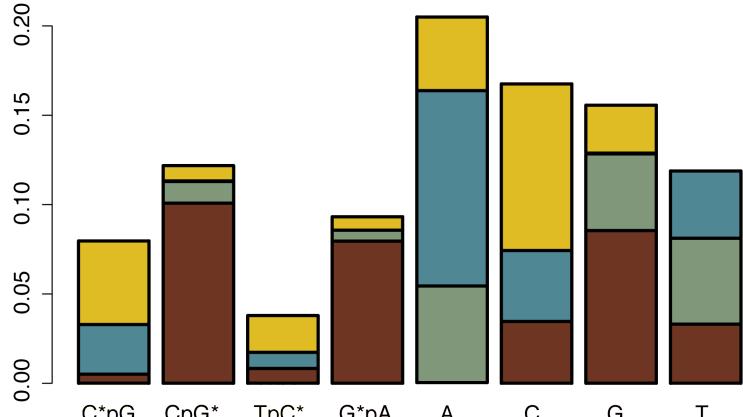
Breast Cancer



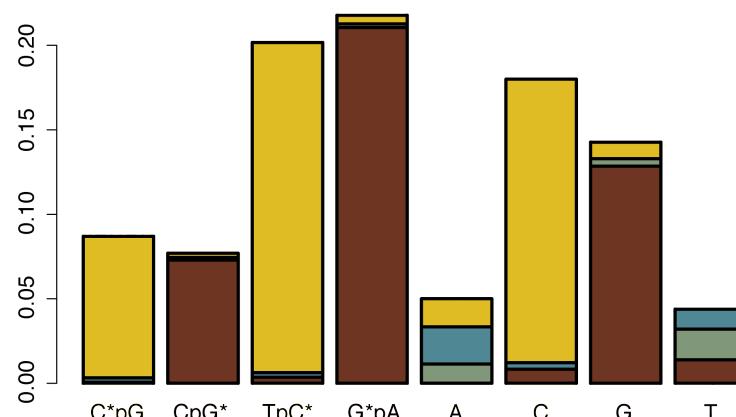
Colon Cancer



Glioblastoma Multiforme



Melanoma



*Infamous
Cancer Genes*

TP53

CDKN2A

KRAS

PIK3CA

PTEN

RB

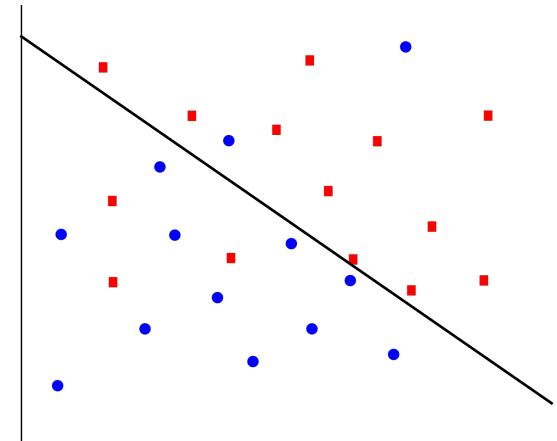
NF1

SMAD4



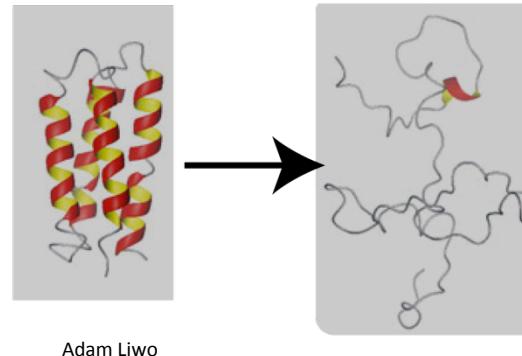
Machine Learning: Supervised Prioritization

- Training Set
 - Positive examples (drivers)
 - Negative examples (passengers)
- Feature set
 - **Informative characteristics of amino acid substitutions**

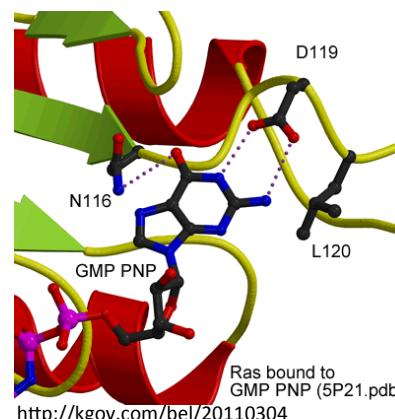


Functional approach: How do amino acid substitutions impact protein activity?

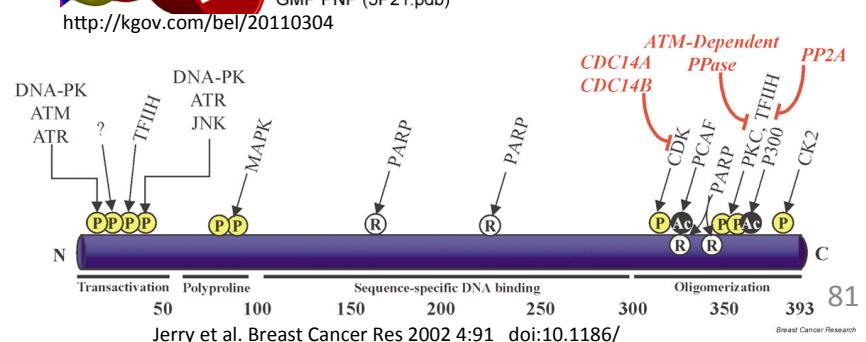
Destabilize the protein



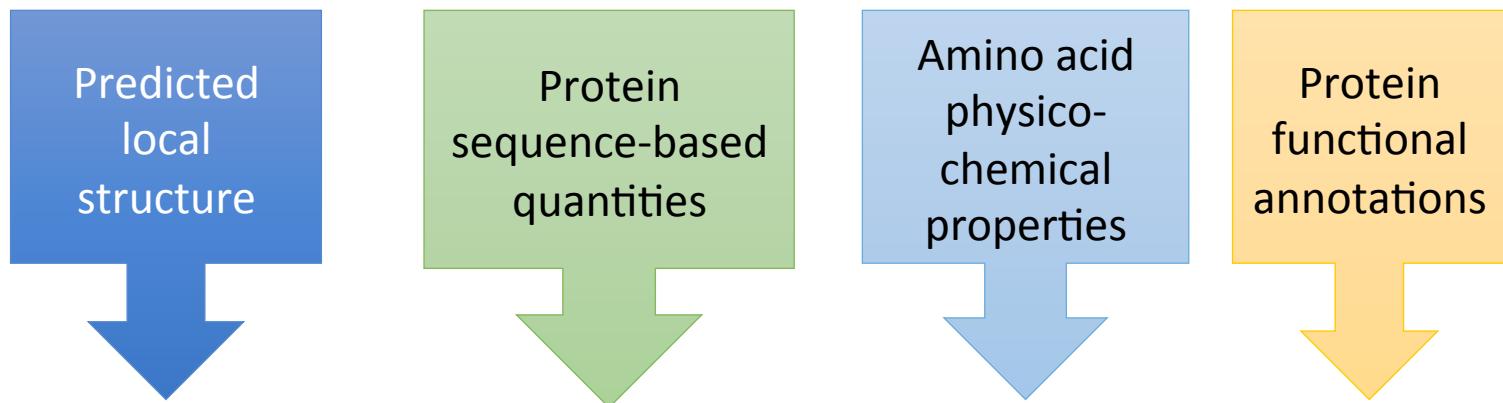
Cause improper localization



Disrupt binding or active sites



Features: Quantifying amino acid substitution impact



Mutation	Residue solvent accessibility	Relative Entropy of an MSA at amino acid position	Charge Change	Binding Site
TP53 S362A	Exposed	4.03	0	Yes
PIK3CA P539R	Buried	4.63	1	No
PIK3CA E545K	Exposed	4.66	2	No

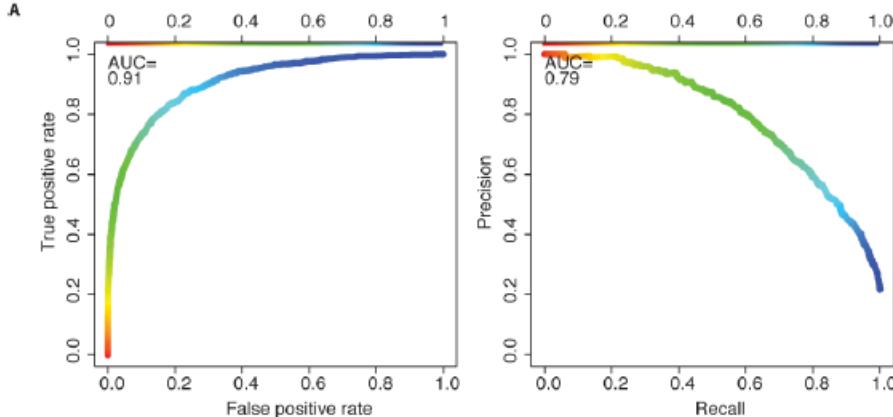
Machine learning algorithm: *Random Forest*

- Ensemble classifier composed of multiple decision trees



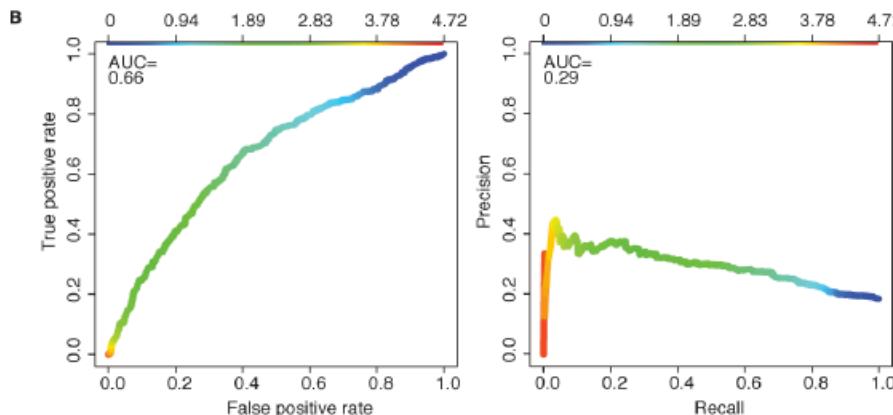
- Each classifier trained on a subset of the training set, using a subset of the features

Comparison to SIFT & PolyPhen

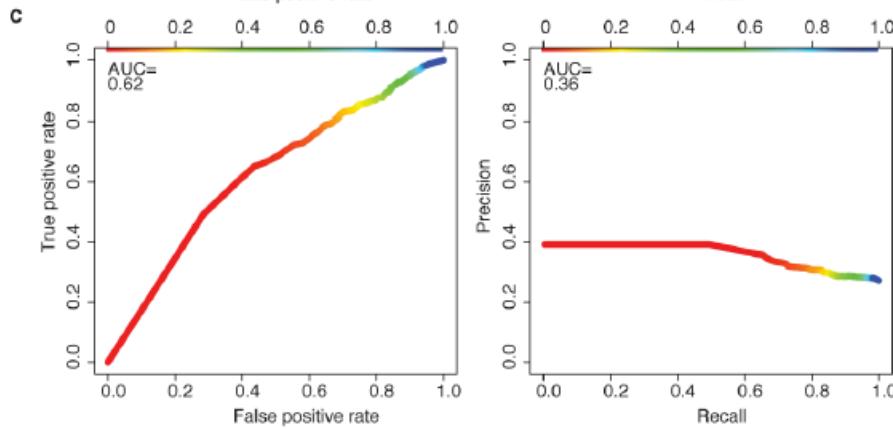


CHASM

Left: ROC Curves
Right: Precision-Recall
Curves



PolyPhen



SIFT

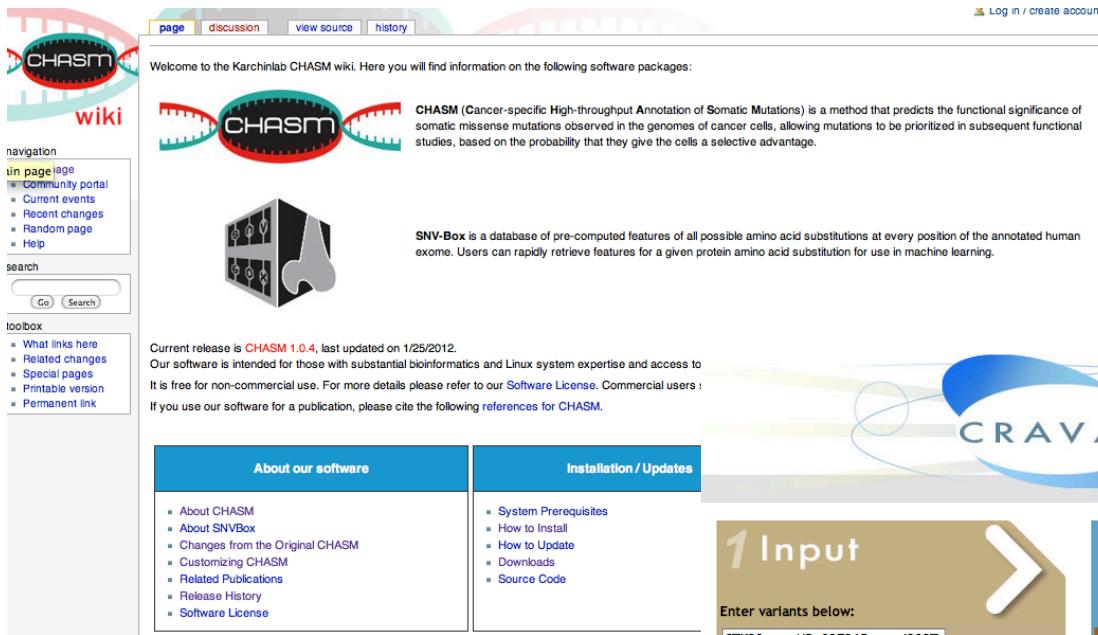
SIFT and PolyPhen score many of the synthetic passengers as likely to alter protein activity

Carter et al Cancer Research 2009



Cancer-specific High-throughput Annotation of Somatic Mutations (CHASM)

Rachel Karchin



Welcome to the Karchinlab CHASM wiki. Here you will find information on the following software packages:

CHASM (Cancer-specific High-throughput Annotation of Somatic Mutations) is a method that predicts the functional significance of somatic missense mutations observed in the genomes of cancer cells, allowing mutations to be prioritized in subsequent functional studies, based on the probability that they give the cells a selective advantage.

SNV-Box is a database of pre-computed features of all possible amino acid substitutions at every position of the annotated human exome. Users can rapidly retrieve features for a given protein amino acid substitution for use in machine learning.

Current release is **CHASM 1.0.4**, last updated on 1/25/2012. Our software is intended for those with substantial bioinformatics and Linux system expertise and access to It is free for non-commercial use. For more details please refer to our [Software License](#). Commercial users : If you use our software for a publication, please cite the following [references for CHASM](#).

About our software

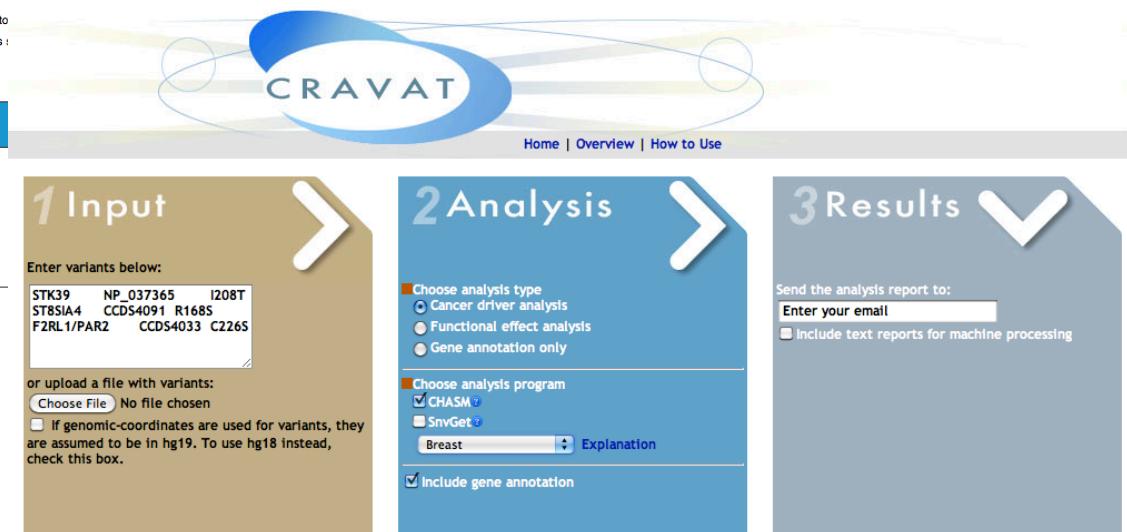
- About CHASM
- About SNVBox
- Changes from the Original CHASM
- Customizing CHASM
- Related Publications
- Release History
- Software License

Installation / Updates

- System Prerequisites
- How to Install
- How to Update
- Downloads
- Source Code

<http://wiki.chasmsoftware.org>

- Random Forest Classifier
- 85 Features
- ~14000 Training examples
- ROC AUC = 0.91



The CRAVAT interface consists of three main steps: Input, Analysis, and Results.

1 Input: Enter variants below:
STK39 NP_037365 I208T
ST8SIA4 CCDS4091 R168S
F2RL1/PAR2 CCDS4033 C226S

or upload a file with variants:
 No file chosen
 If genomic-coordinates are used for variants, they are assumed to be in hg19. To use hg18 instead, check this box.

2 Analysis: Choose analysis type
 Cancer driver analysis
 Functional effect analysis
 Gene annotation only

Choose analysis program
 CHASM
 SnvGet
Breast Explanation
 Include gene annotation

3 Results: Send the analysis report to:

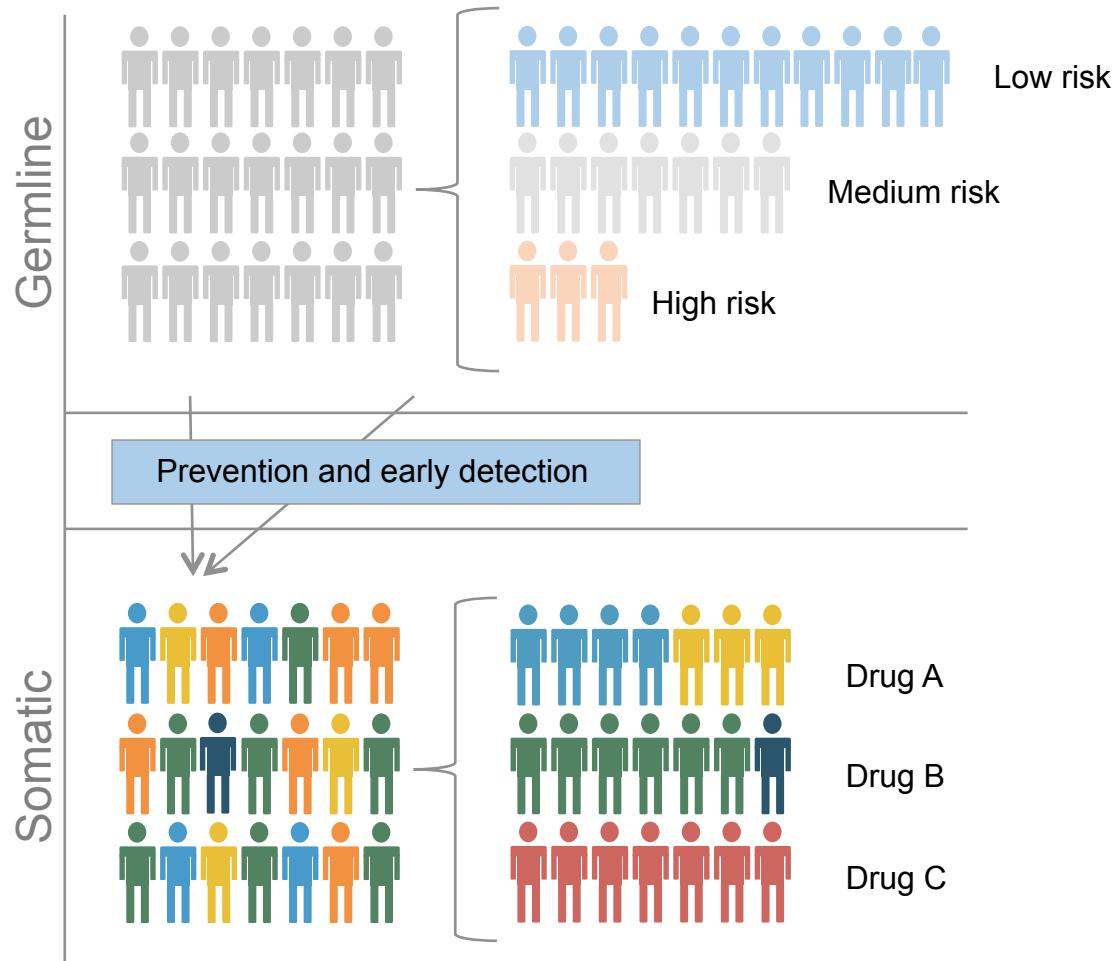
 Include text reports for machine processing

<http://www.cravat.us>

SUBMIT

Precision Cancer Medicine

- Identify individuals at risk
- Preventative measures and screening for early detection
- Patient stratification for prognostic or treatment purposes



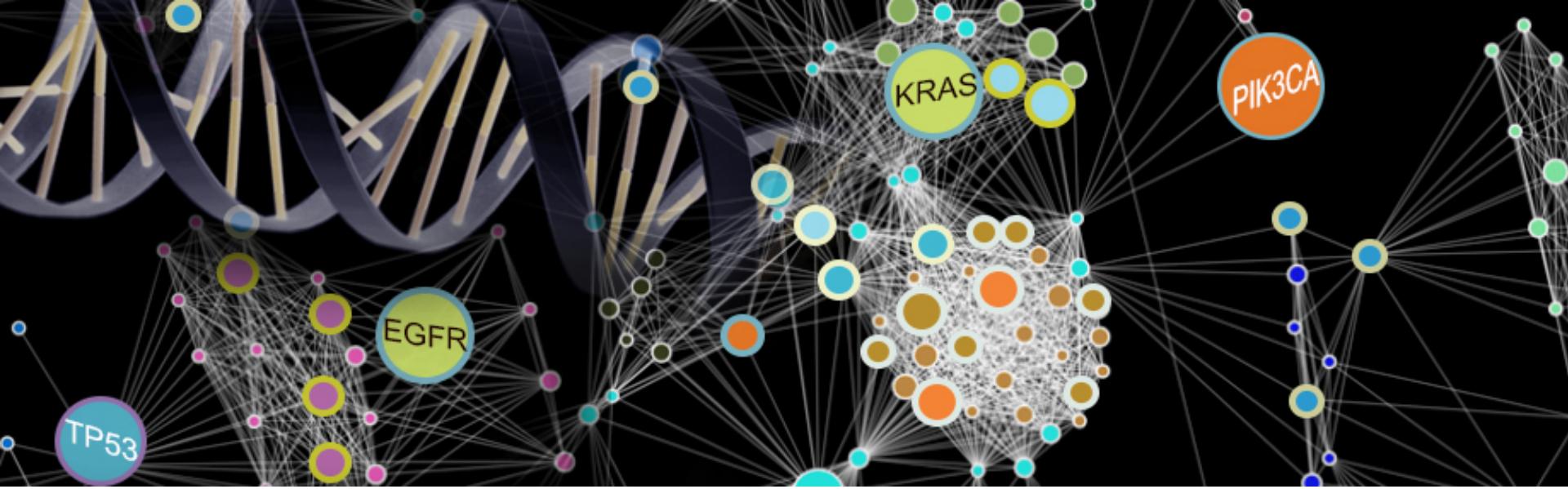
Conclusions

Mutation burden is more than just the number of mutations in a tumor

The majority of mutations in cancer are passengers – it is important to know which are drivers

Cancer evolution generates heterogeneity in many ways

Quantifying these factors is necessary to develop effective treatment strategies



THANK YOU!