

Reasoning about biology with data-driven approaches

Howes Scholar Presentation

Adam Riesselman

Reasoning about biology with data-driven approaches

Howes Scholar Presentation

Adam Riesselman

DOE CSGF Fellow: Harvard University, 2014-2018

Practicum: Lawrence Berkeley National Lab
(Joint Genome Institute), 2016

Understanding how **genotype** affects **phenotype** is a critical **challenge**.

T
CATAACCA↑GTACA
C

Understanding how **genotype** affects **phenotype** is a critical **challenge**.

T
CATAACCA↑GTACA
C

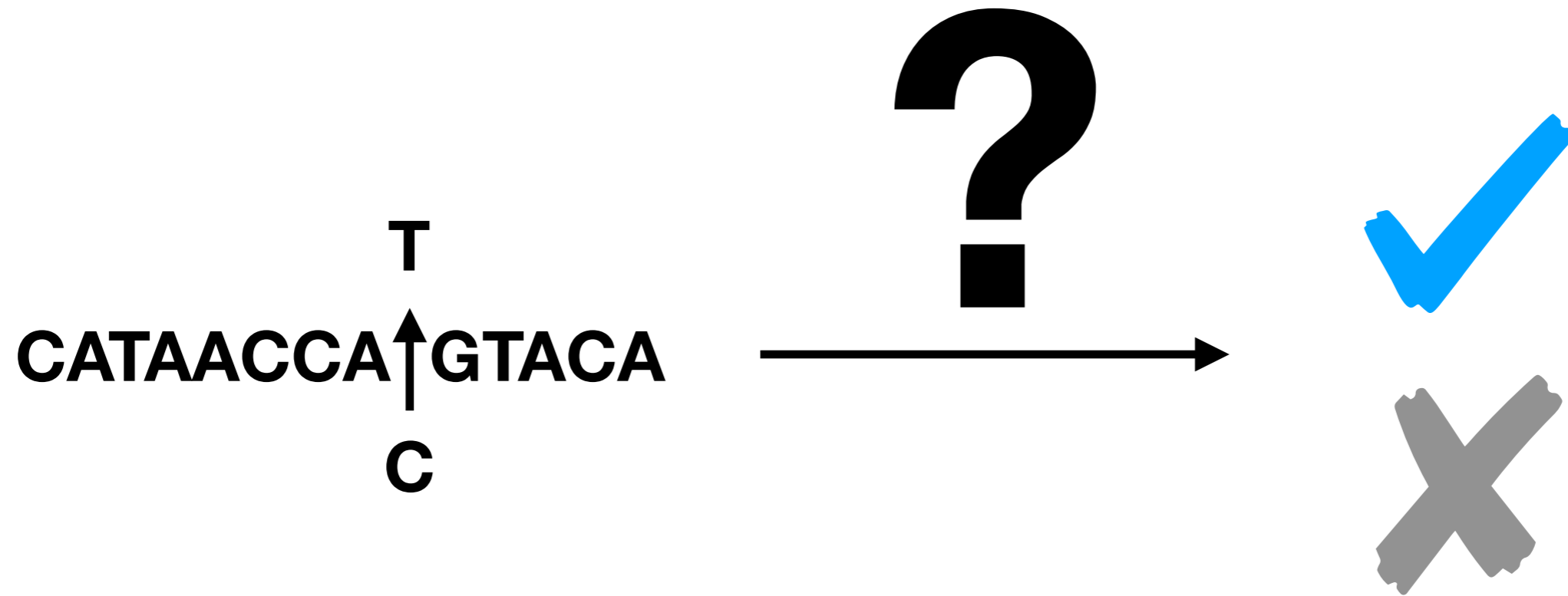


Understanding how **genotype** affects **phenotype** is a critical **challenge**.

T
CATAACCA↑GTACA
C

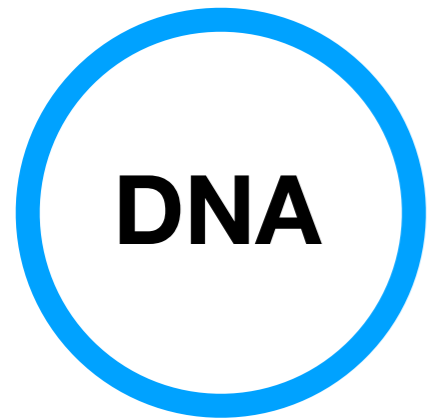


Understanding how **genotype** affects **phenotype** is a critical **challenge**.

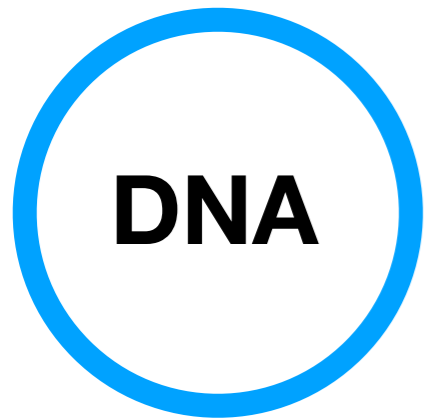


The **Central Dogma** of **Biology**

The **Central Dogma** of **Biology**



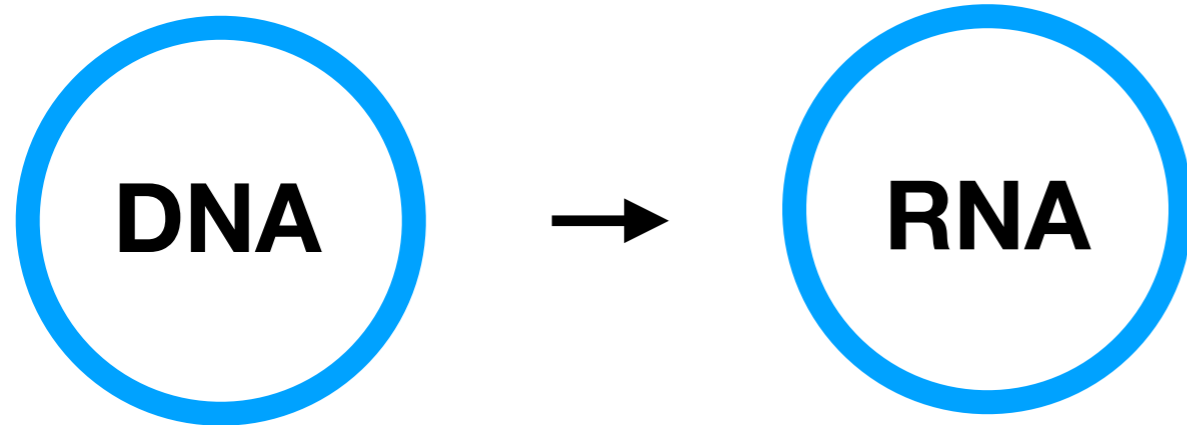
The Central Dogma of Biology



ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTCCGCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTA ACTCGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATAACCAACGACGAG...

4 different
bases

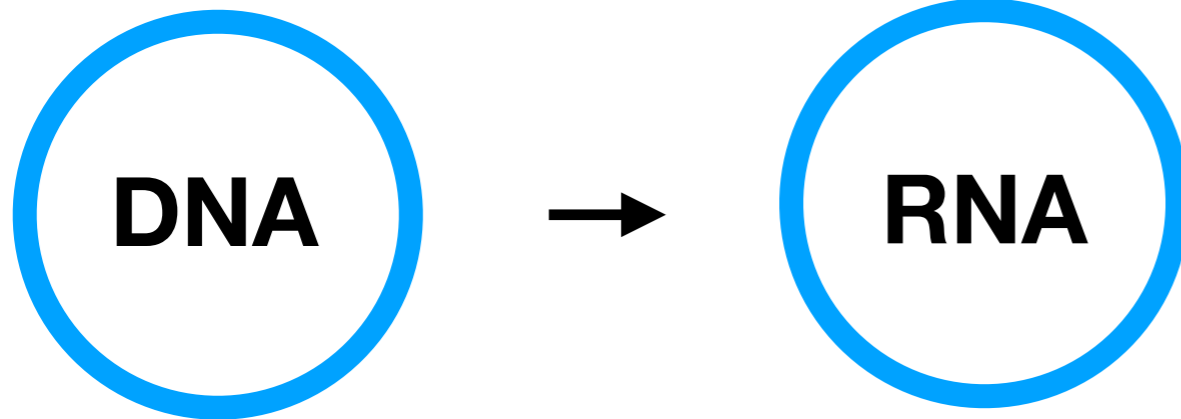
The Central Dogma of Biology



ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTCCGCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTA ACTCGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATAACCAACGACGAG...

4 different
bases

The Central Dogma of Biology

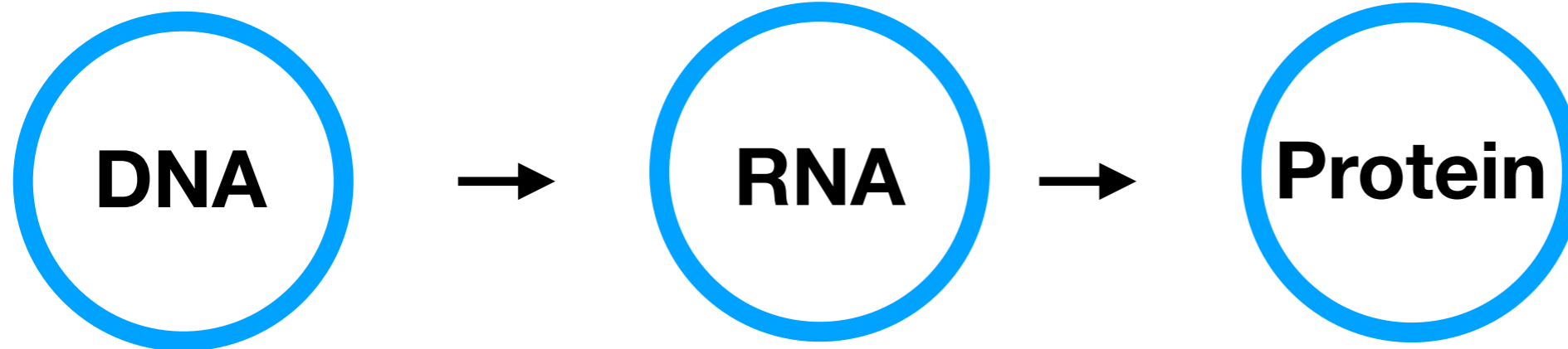


ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTCCGCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATAACCAACGACGAG...

AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCGAAG
AACGUUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUCACCAGUCACAG
AAAAGCAUCUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUAACA
CUGCGGCCAACUUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAACGACGAG...

4 different
bases

The Central Dogma of Biology

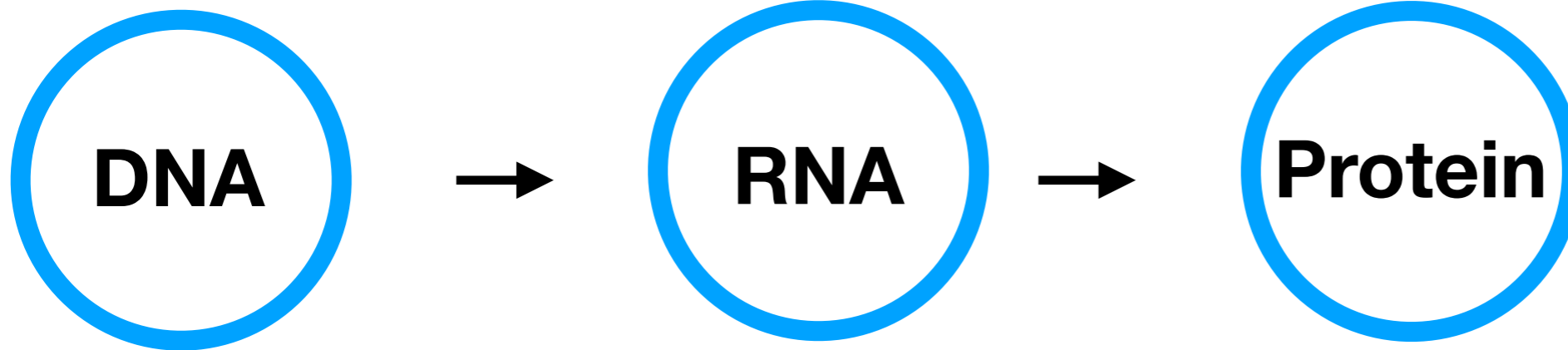


ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTCCGCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATAACCAACGACGAG...

AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCGAAG
AACGUUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUCACCAGUCACAG
AAAAGCAUCUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUAACA
CUGCGGCCAACUUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAACGACGAG...

4 different
bases

The Central Dogma of Biology



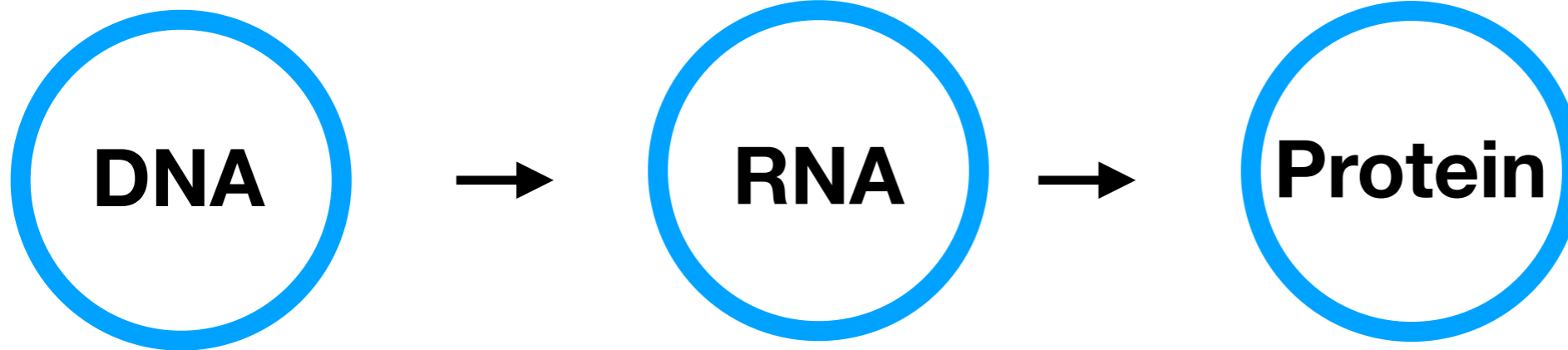
ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTGCGCCG
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATAACCAACGACGAG...

AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCGAAG
AACGUUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUCACCAGUCACAG
AAAAGCAUCUUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUAACA
CUGCGGCCAACUUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAACGACGAG...

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAFLHNMGDHVTRL
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSAIPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

20 different
amino acids

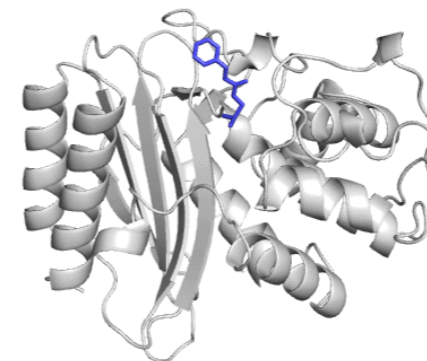
The Central Dogma of Biology



ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTGCGCCG
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATACCAAACGACGAG . . .

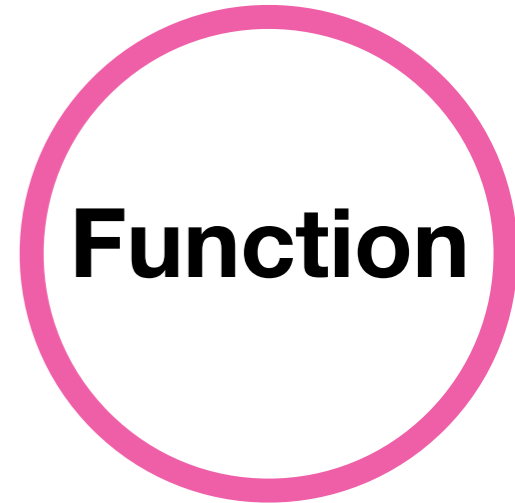
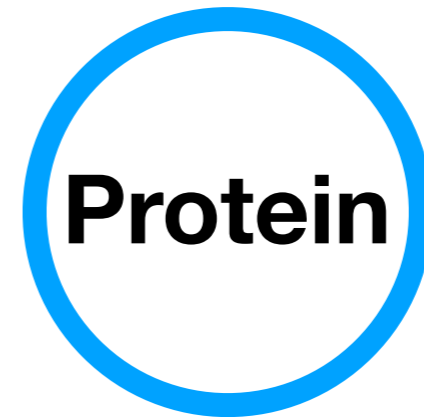
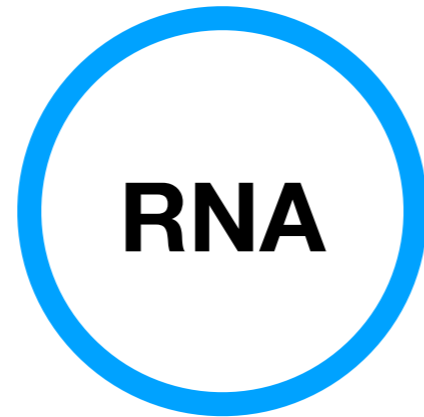
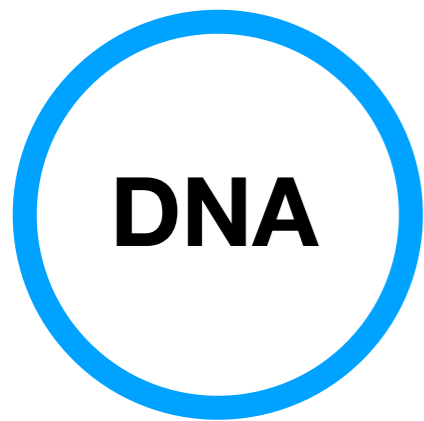
AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCGAAG
AACGUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUCACCAGUCACAG
AAAAGCAUCUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUACA
CUGCGGCCAACUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAAACGACGAG . . .

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGDHSVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLR SALPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



20 different
amino acids

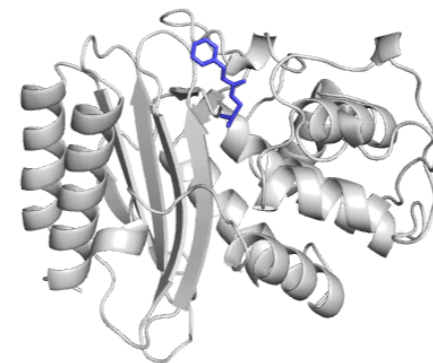
The Central Dogma of Biology



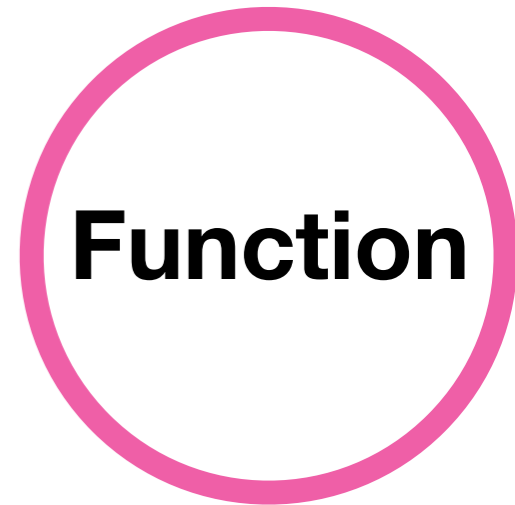
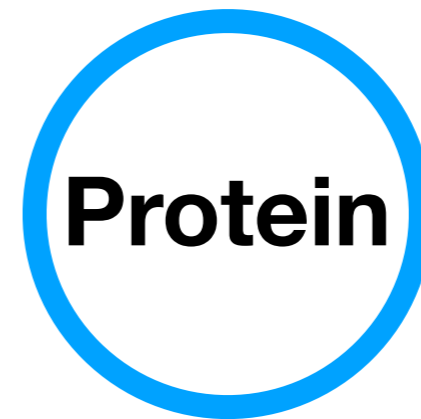
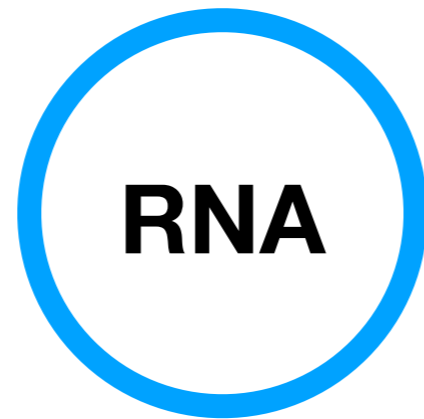
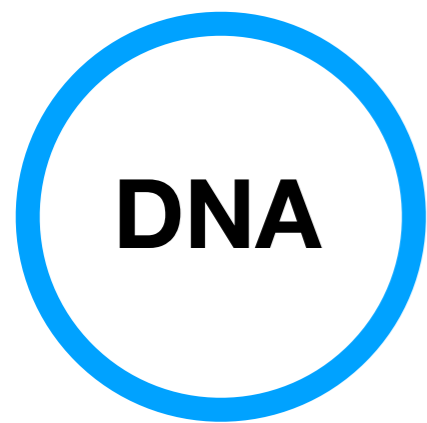
ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTGCGCCG
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCGATCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATACCAAACGACGAG...

AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCGAAG
AACGUUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUCACCAGUCACAG
AAAAGCAUCUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUAACA
CUGCGGCCAACUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAAACGACGAG...

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGDHSVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLTTGELLTLASRQQLID
WMEADKVVAGPLLRSALPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



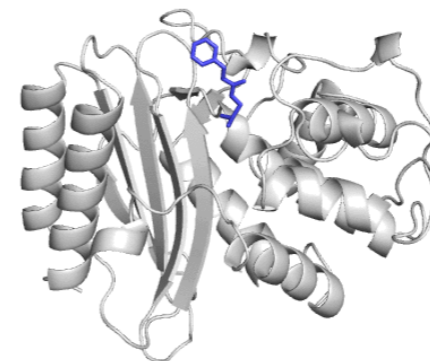
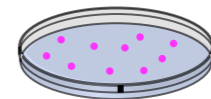
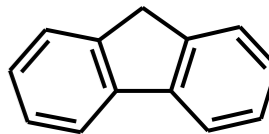
The Central Dogma of Biology



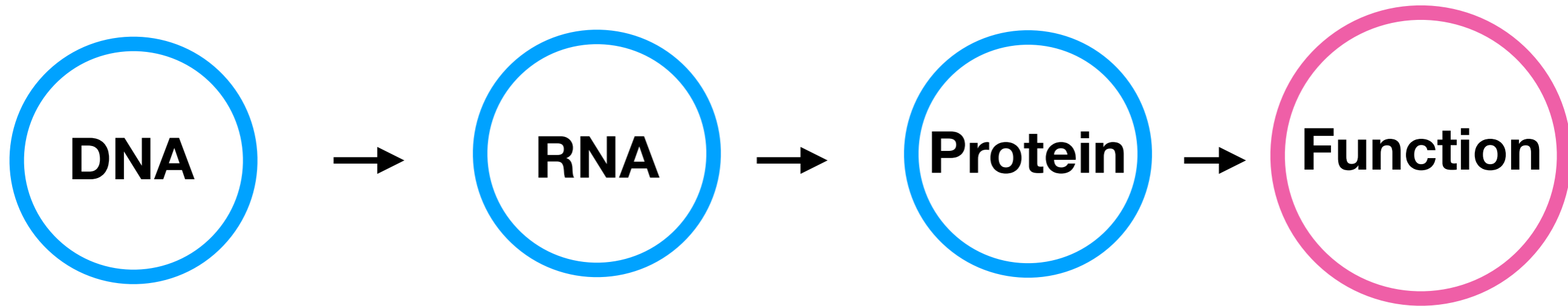
ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTTGACGCCG
GGCAAGAGCAACTCGGTGCGCCG
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTAACCTGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATACCAAACGACGAG...

AUGAGUAUUCAACAUUUCGUGU
CGCCCUUAUUCUUUUUUGCGG
CAUUUUGCCUUCUGUUUUUGCU
CACCCAGAAACGCUGGUGAAAGU
AAAAGAUGCUGAAGAUCAGUUGG
GUGCACGAGUGGGUACAUCGAA
CUGGAUCUCAACAGCGGUAAGAU
CCUUGAGAGUUUCGCCCCGAAG
AACGUUUUCCAAUGAUGAGCACU
UUUAAAGUUCUGCUAUGUGGCGC
GGUAUUAUCCCGUGUUGACGCCG
GGCAAGAGCAACUCGGUCGCCG
AUACACUAUUCUCAGAAUGACUU
GGUUGAGUACUACCAGUCACAG
AAAAGCAUCUACGGGAUGGCAUG
ACAGUAAGAGAAUUAUGCAGUGC
UGCCAUAACCAUGAGUGAUACA
CUGCGGCCAACUACUUCUGACA
ACGAUCGGAGGACCGAAGGAGCU
AACCGCUUUUUUGCACAACAUGG
GGGAUCAUGUAACUCGCCUUGAU
CGUUGGGAACCGGAGCUGAAUGA
AGCCAUAACCAAACGACGAG...

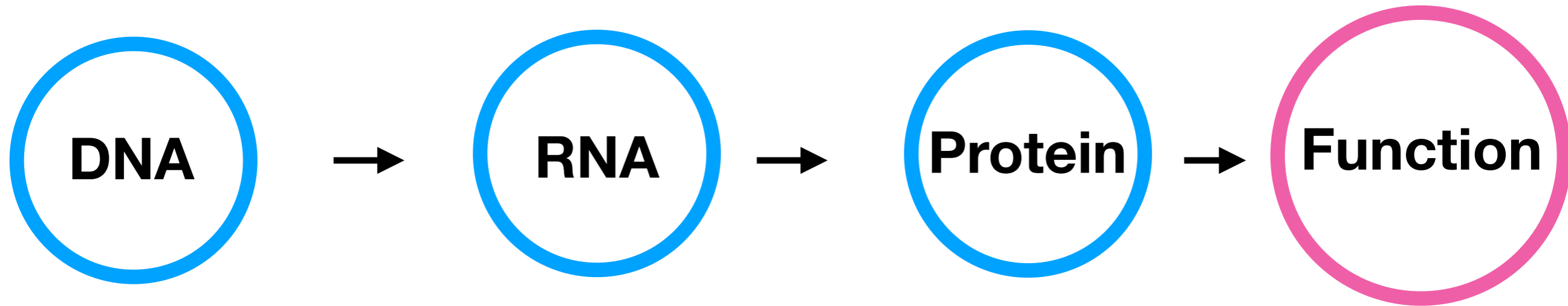
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGDHSVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLR SALPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



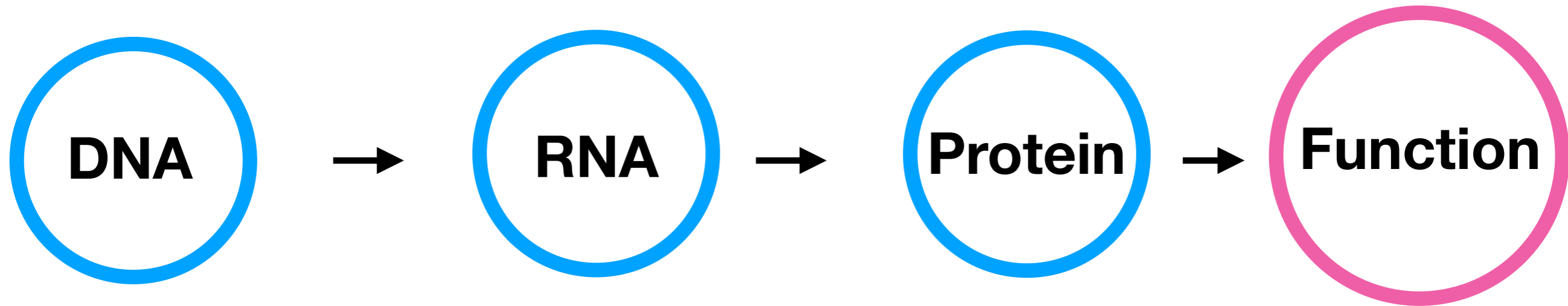
The Central Dogma of Biology



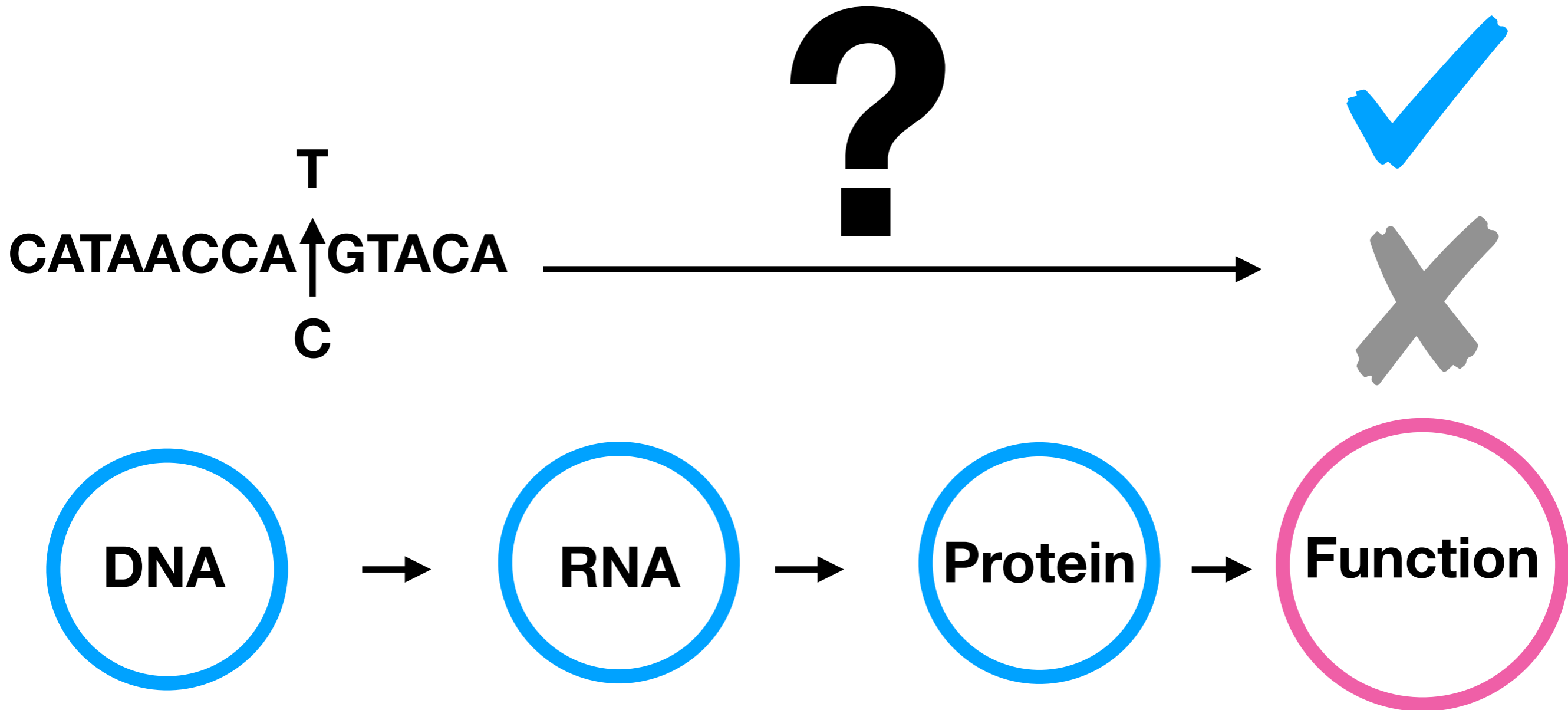
Mutations alter biological **function**



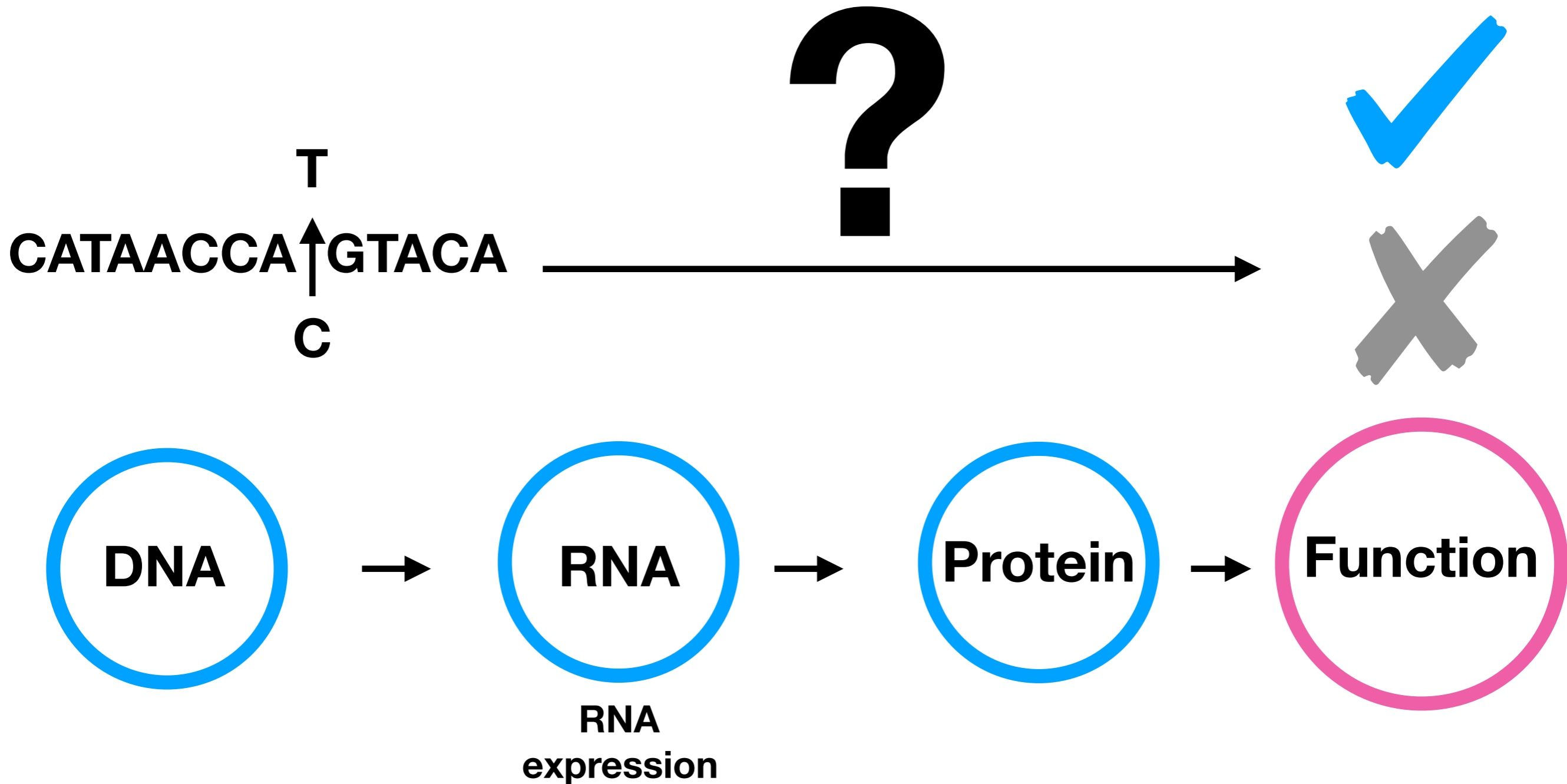
Mutations alter biological **function**



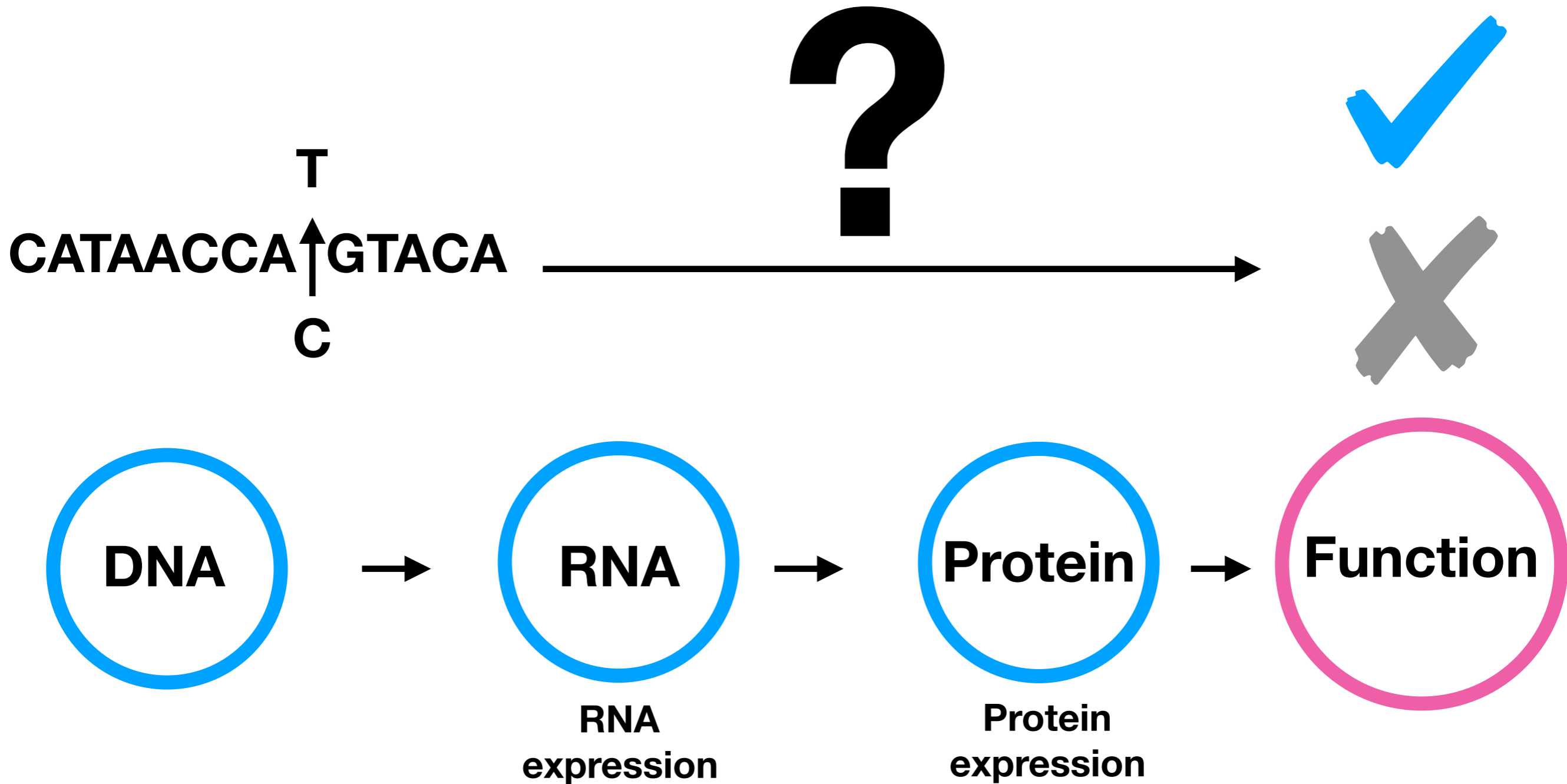
Mutations alter biological function



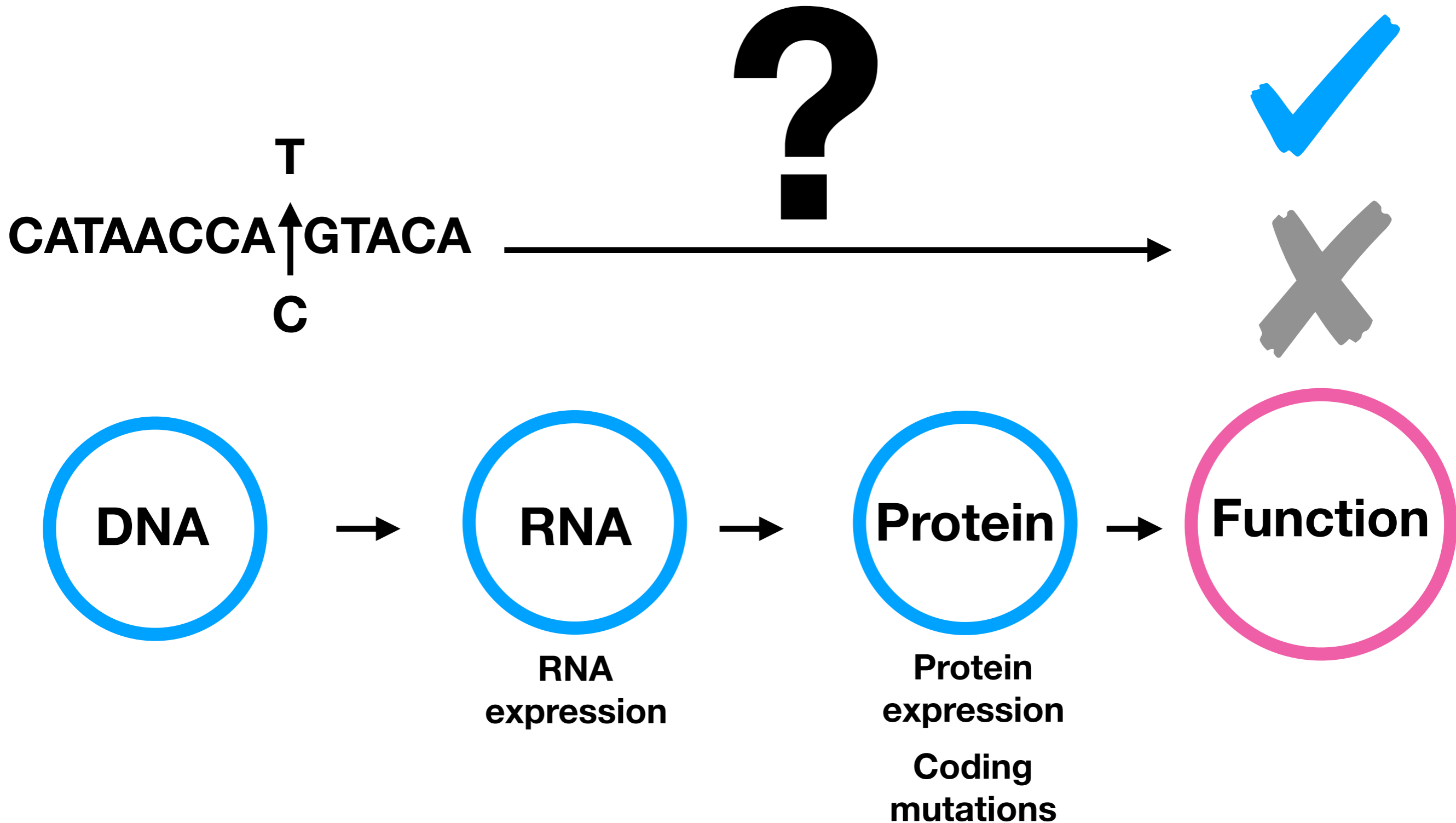
Mutations alter biological function



Mutations alter biological function



Mutations alter biological function



**Mutation effect prediction
is important**

Mutation effect prediction is important

Understanding
disease



“Does this mutation
cause cancer?”

Mutation effect prediction is important

Understanding
disease

Biomedicine



“Does this mutation
cause cancer?”

“Is this antibody
stable in a patient?”

Mutation effect prediction is important

Understanding
disease

Biomedicine

Bioengineering

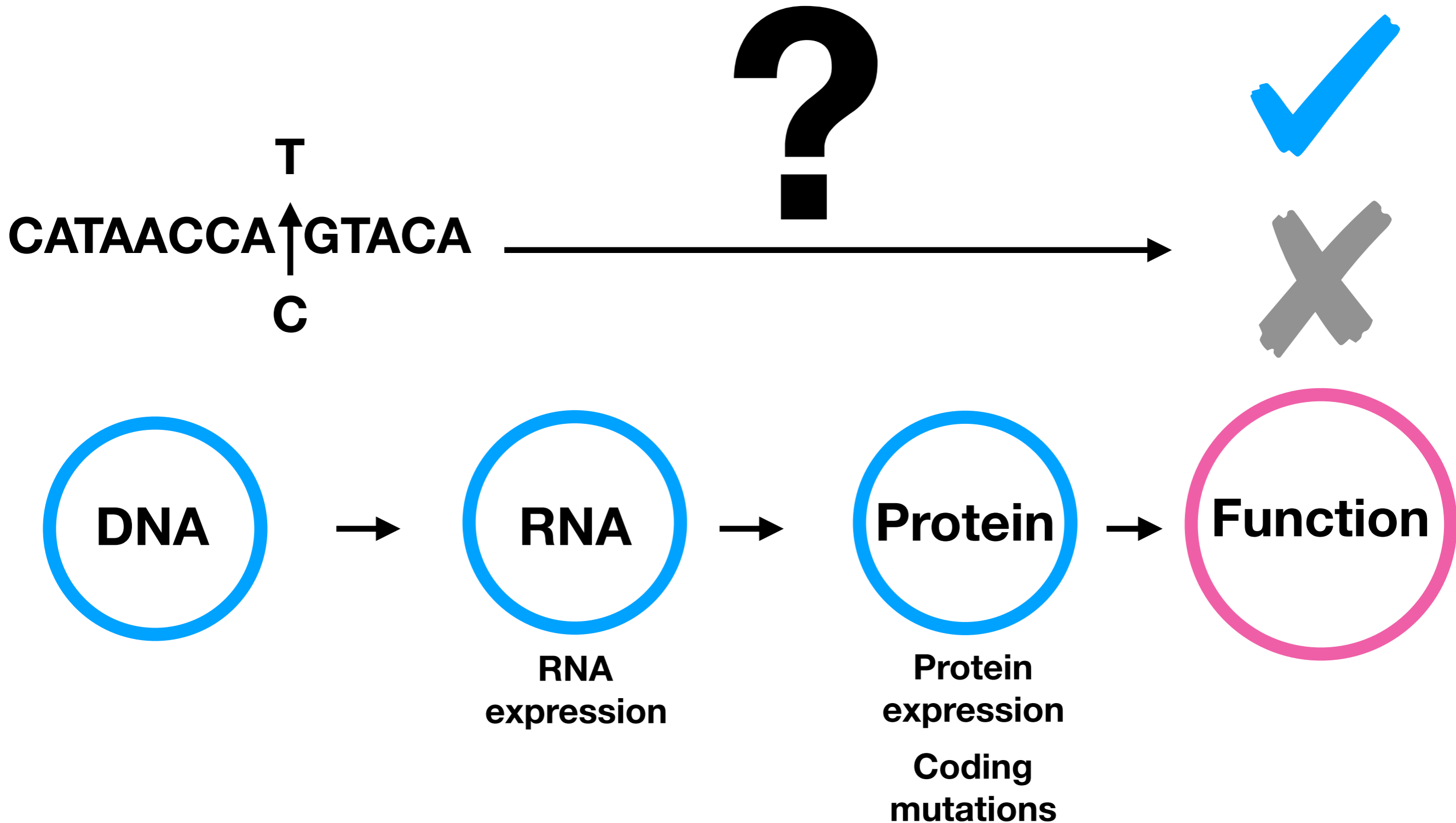


“Does this mutation
cause cancer?”

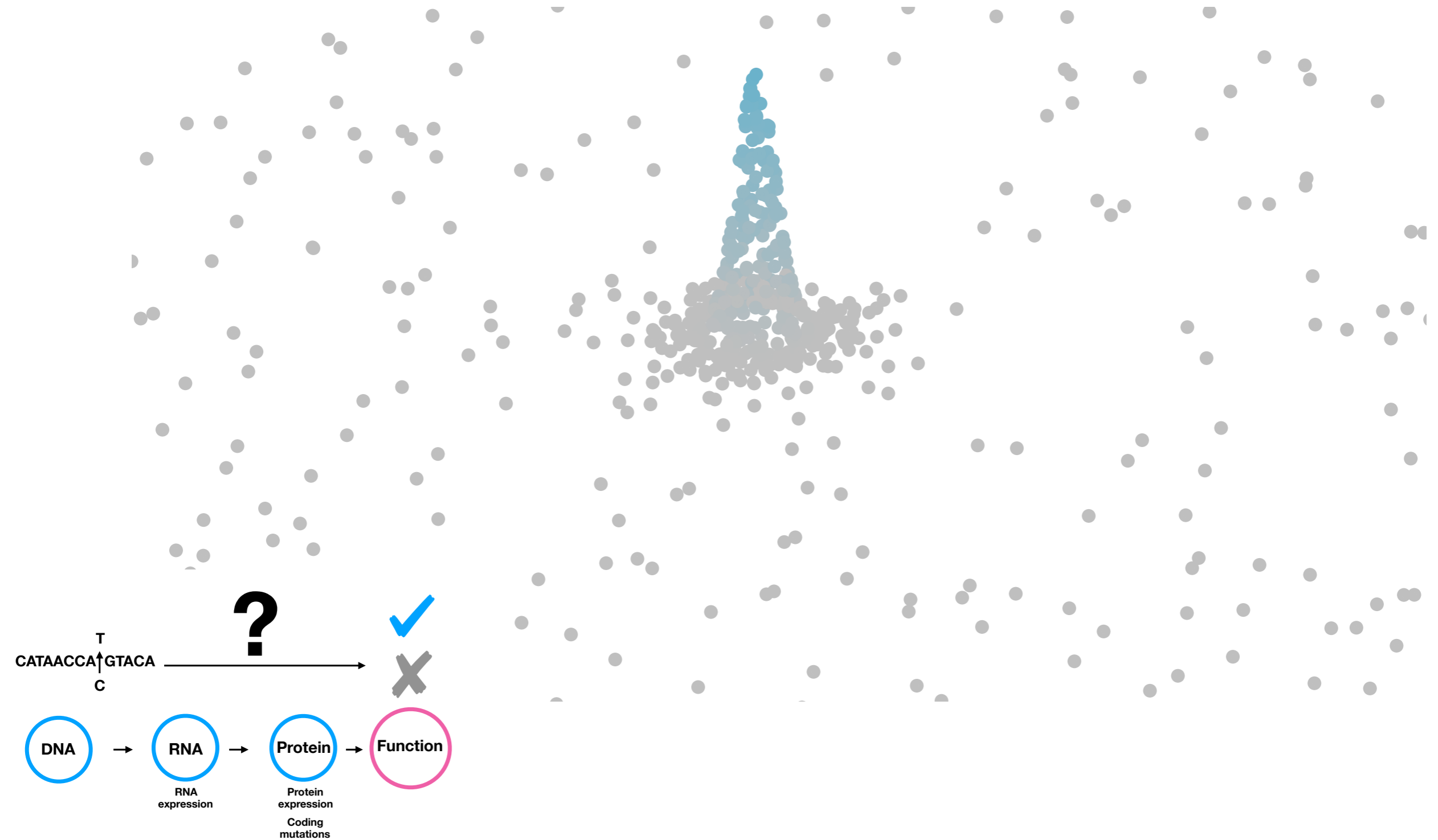
“Is this antibody
stable in a patient?”

“Can this microbe be
used to create
vitamins?”

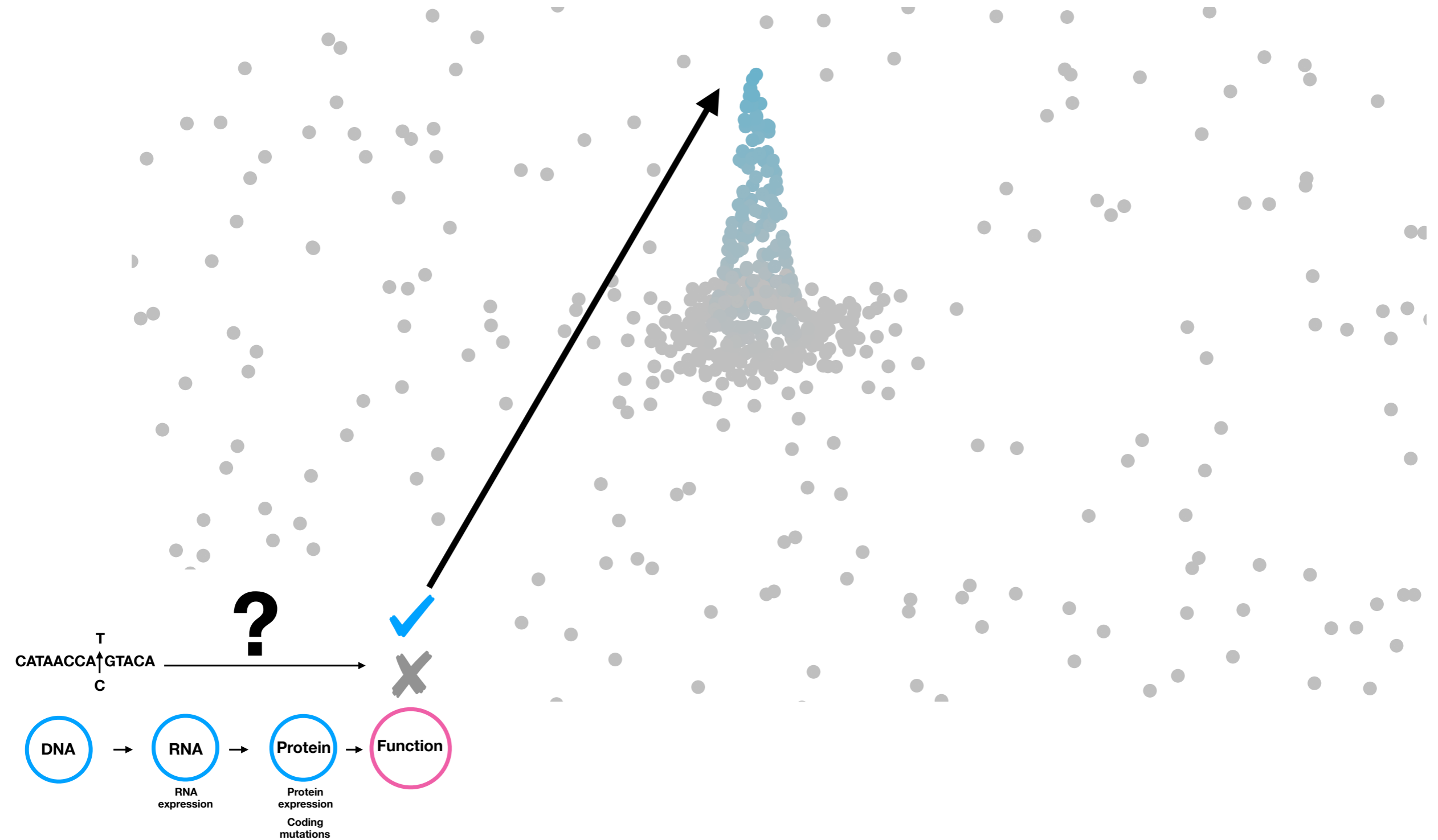
Mutations alter biological function



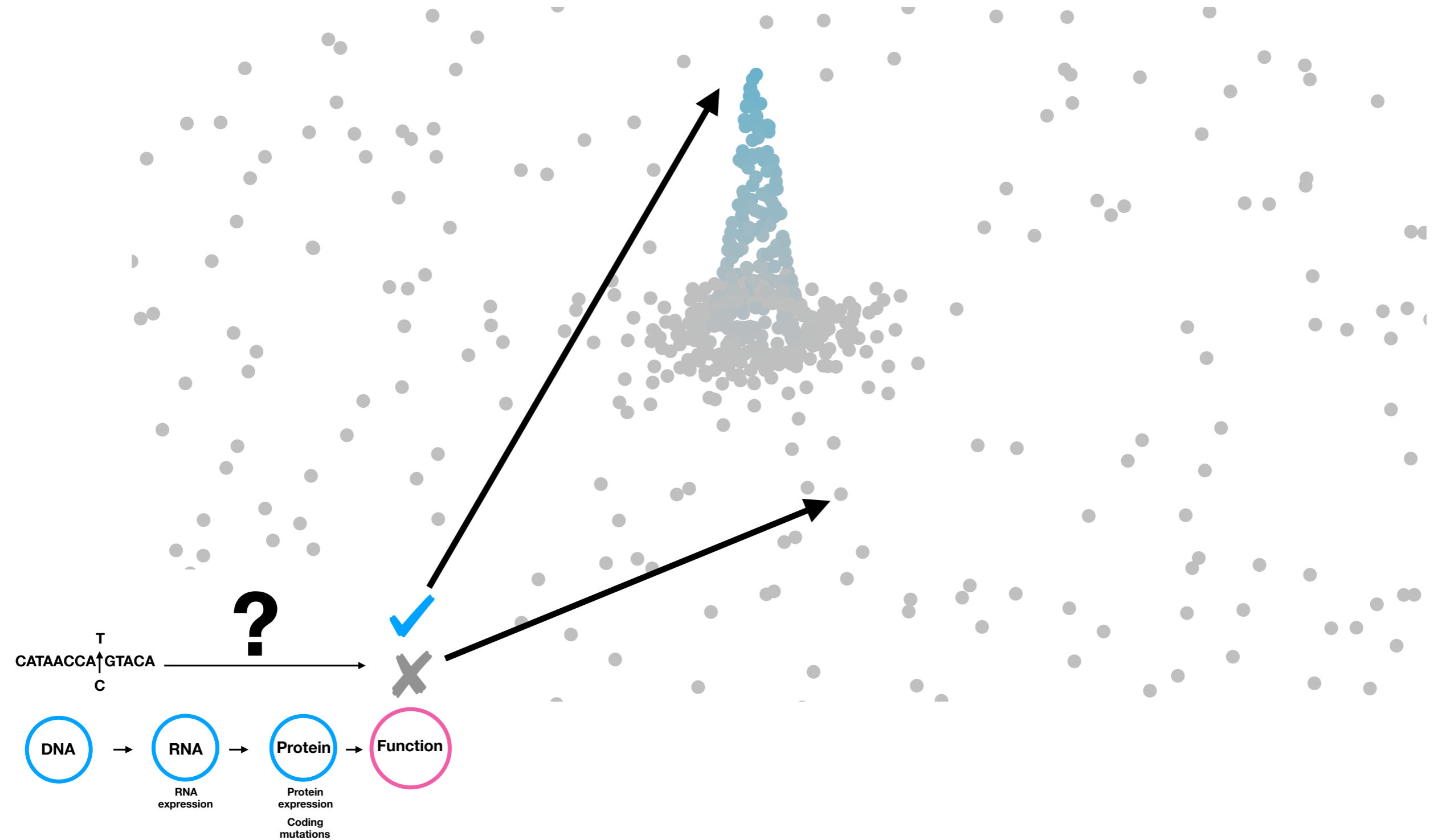
Mutations alter biological function



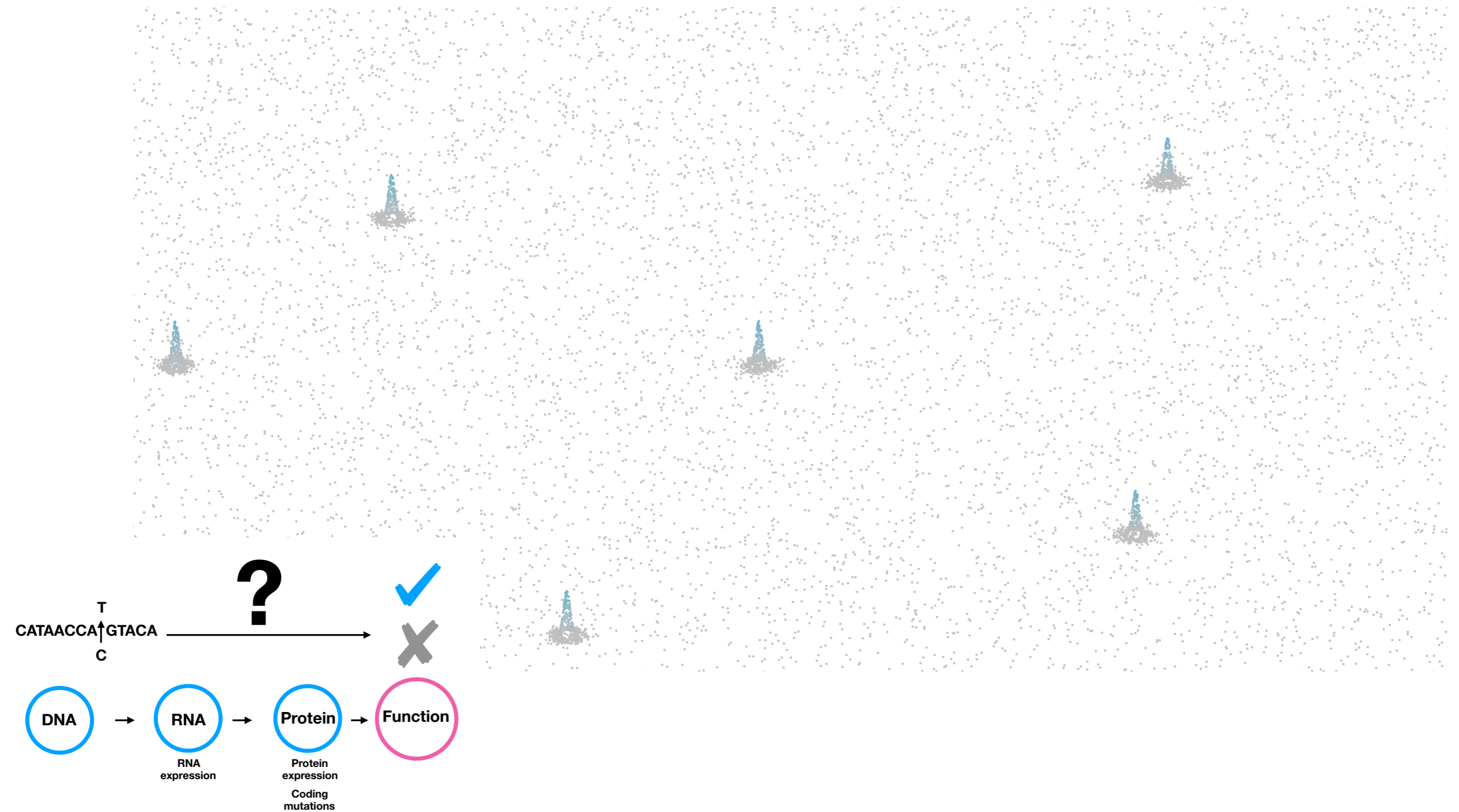
Mutations alter biological function



Mutations alter biological function



Mutations alter biological function



Mutations alter biological function

CATAACCA↑GTACA
 T
 C



Mutation effect prediction is hard

CATAACCA↑GTACA
 T
 C



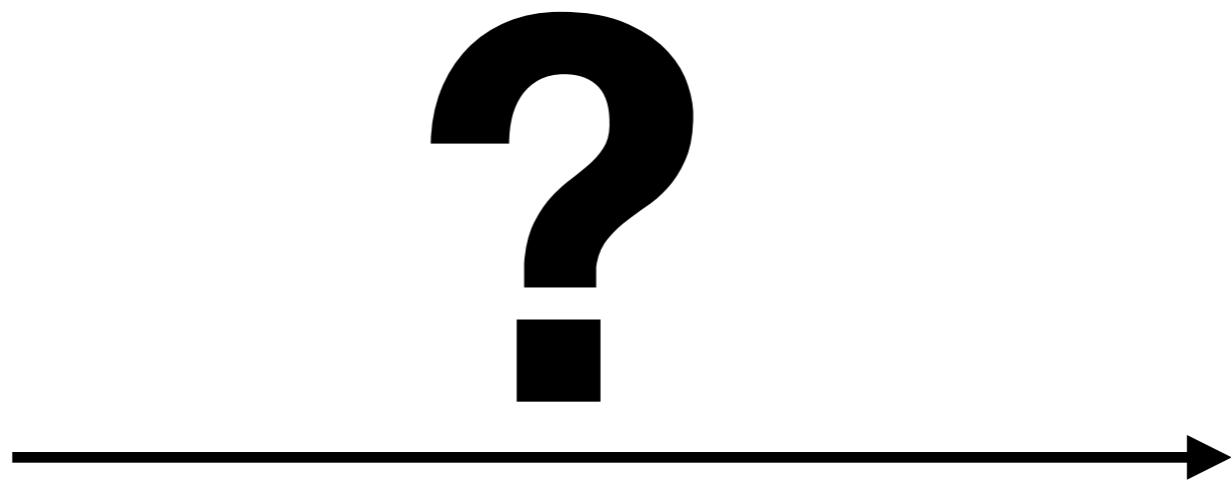
?



Mutation effect prediction is hard

Sparsely sampled

CATAACCA↑GTACA
 T
 C

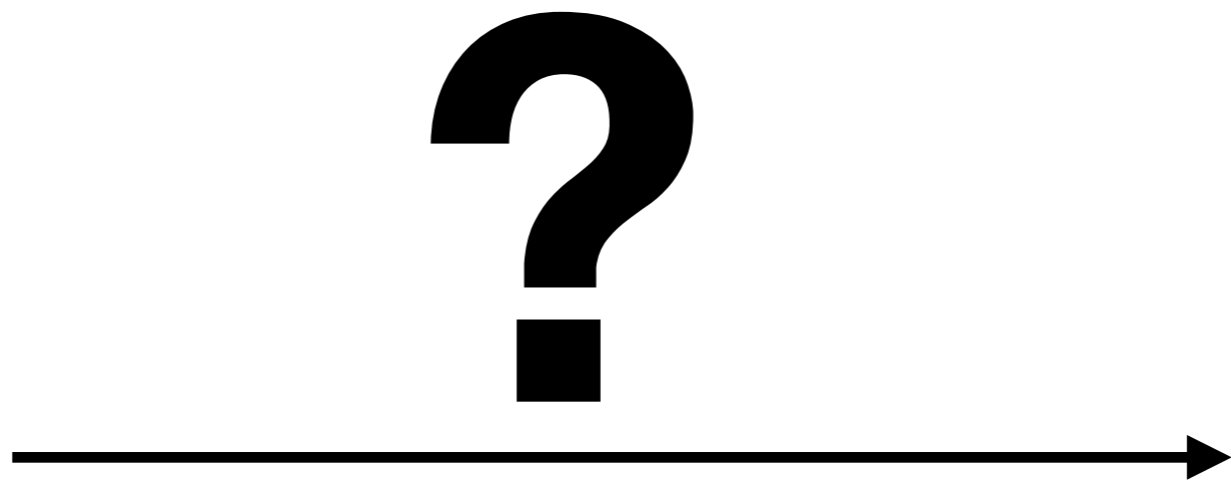


Mutation effect prediction is hard

Sparsely sampled

Nonlinear interactions

CATAACCA↑GTACA
 T
 C



Mutation effect prediction is hard

Sparsely sampled

Noisy

Nonlinear interactions

T
CATAACCA↑GTACA
C



?



Mutation effect prediction is hard

Sparsely sampled

Noisy

Nonlinear interactions

Confounders

T
CATAACCA↑GTACA
C



?



Mutation effect prediction is hard

Sparsely sampled

Noisy

Effect not measured

Nonlinear interactions

Confounders

T
CATAACCA↑GTACA
C



Part I:

Genotype -> Phenotype
in proteins

Part I: Genotype -> Phenotype in proteins

DNA

ATGAGTATTCAACATTTCCGTGT
CGCCCTTATTCCCTTTTTTGC GG
CATTTTGCCTTCCTGTTTTTGCT
CACCCAGAAACGCTGGTGAAAGT
AAAAGATGCTGAAGATCAGTTGG
GTGCACGAGTGGGTTACATCGAA
CTGGATCTCAACAGCGGTAAGAT
CCTTGAGAGTTTTTCGCCCCGAAG
AACGTTTTCCAATGATGAGCACT
TTTAAAGTTCTGCTATGTGGCGC
GGTATTATCCCGTGTGACGCCG
GGCAAGAGCAACTCGGTCGCCGC
ATACACTATTCTCAGAATGACTT
GGTTGAGTACTCACCAGTCACAG
AAAAGCATCTTACGGATGGCATG
ACAGTAAGAGAATTATGCAGTGC
TGCCATAACCATGAGTGATAACA
CTGCGGCCAACTTACTTCTGACA
ACGATCGGAGGACCGAAGGAGCT
AACCGCTTTTTTGCACAACATGG
GGGATCATGTA ACTCGCCTTGAT
CGTTGGGAACCGGAGCTGAATGA
AGCCATACCAAACGACGAG...

RNA

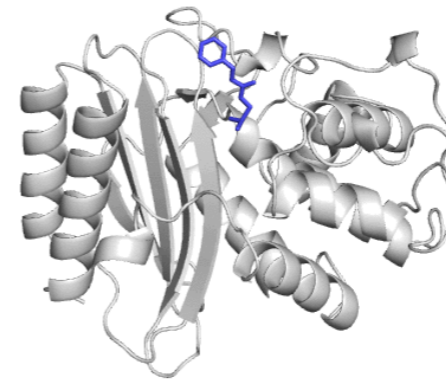


Protein

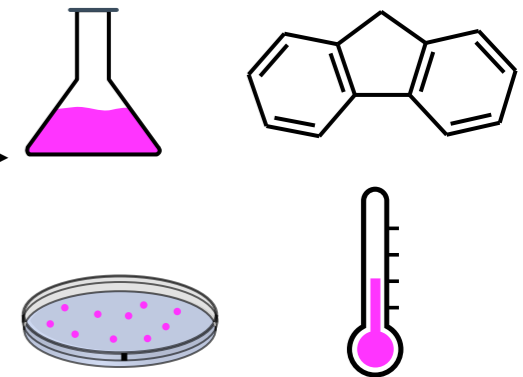
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKLLT GELLTLASRQQLID
WMEADK VAGPLLRSALPAGWFIA
DKSGAGERGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



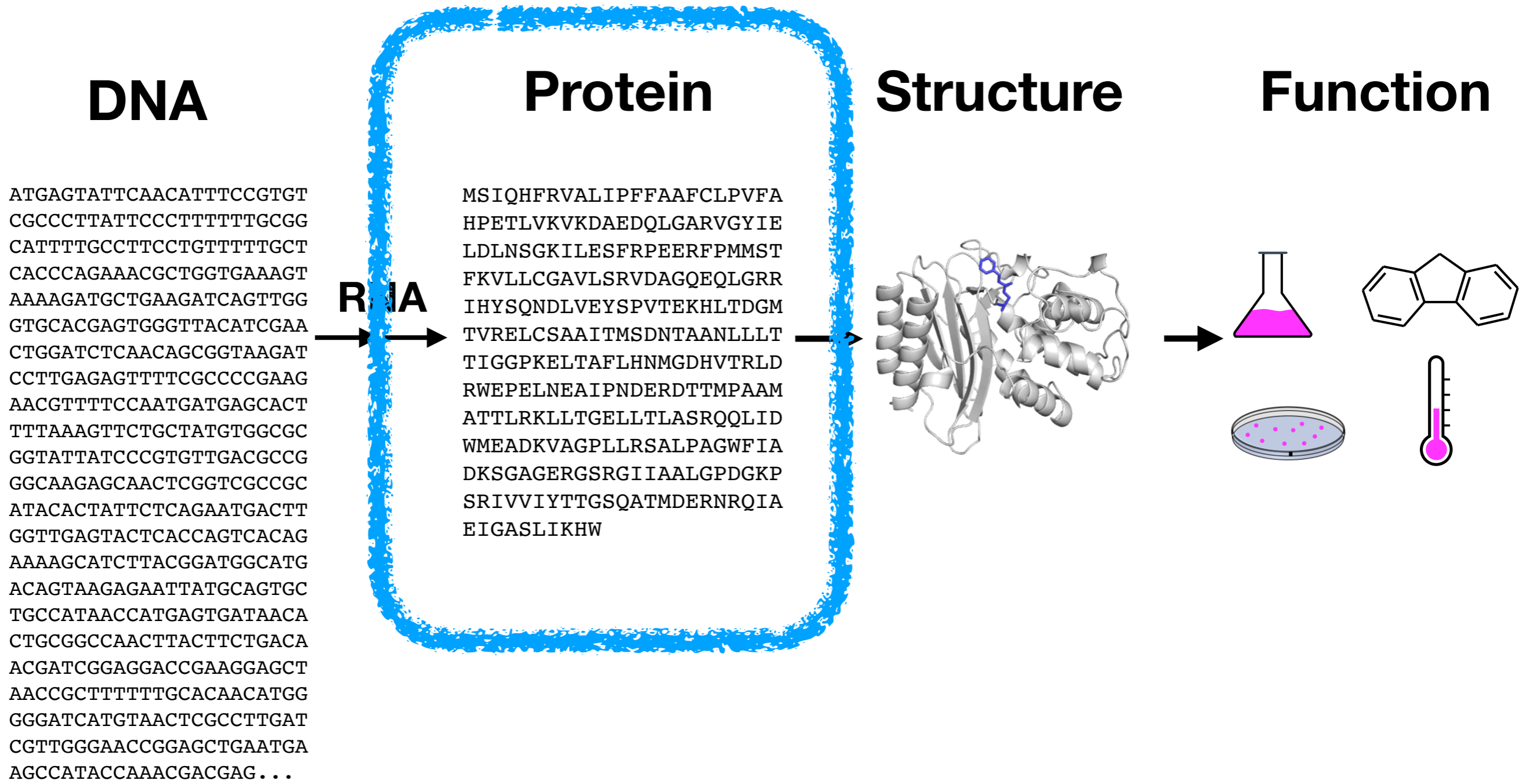
Structure



Function



Part I: Genotype -> Phenotype in proteins



Mutations impact protein **function**

Mutations impact protein function

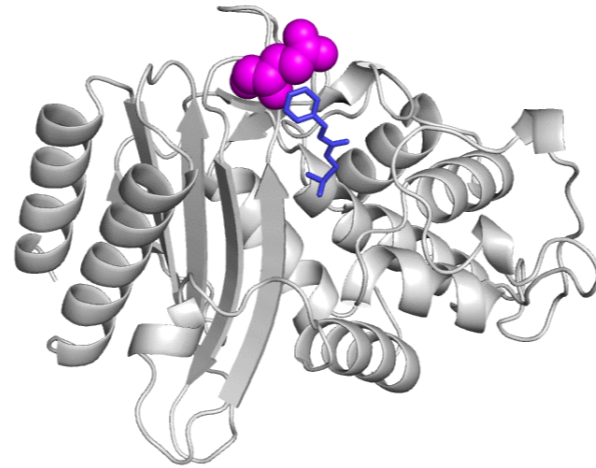
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAGE^RGRSGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

Sequence

Mutations impact protein function

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAFLHNMGDHVTRLR
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAGE^RRGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



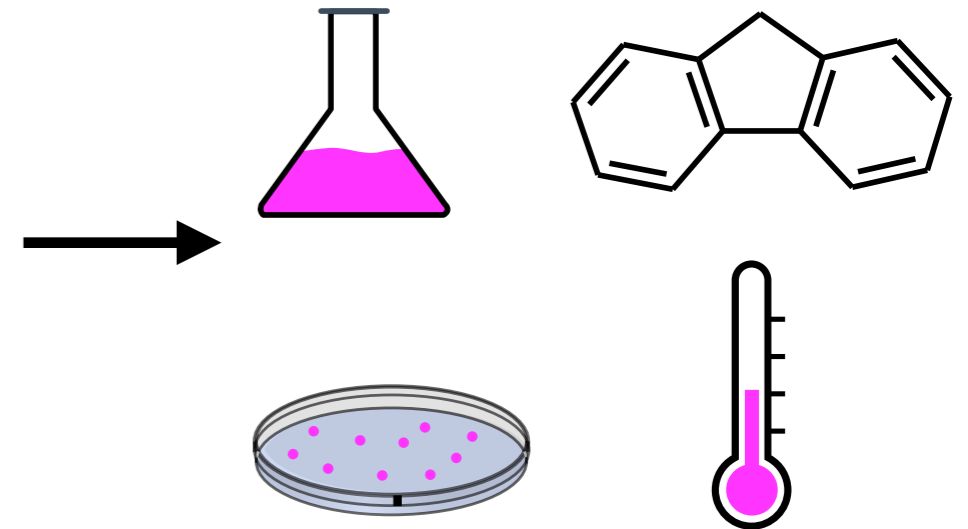
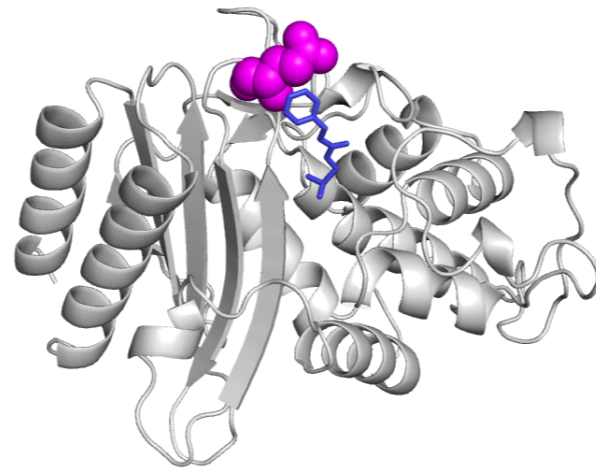
Sequence

Structure

Mutations impact protein function

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



Sequence

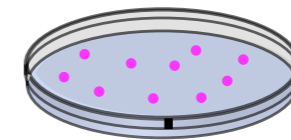
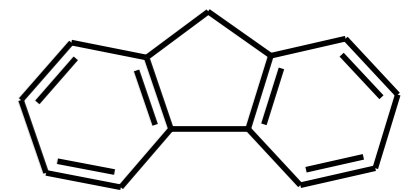
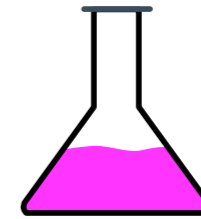
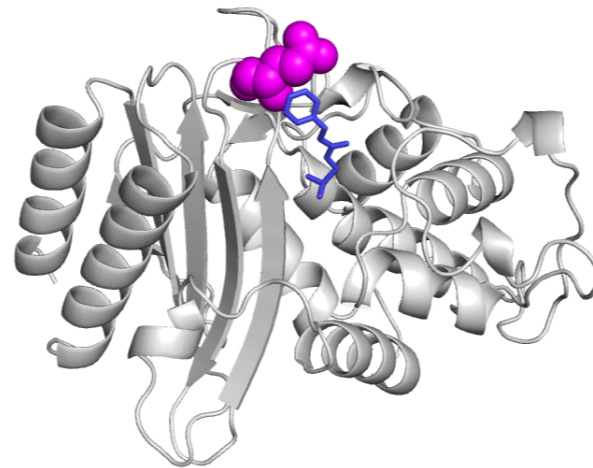
Structure

Function

Mutations impact protein function

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **E** RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **L** RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

Sequence

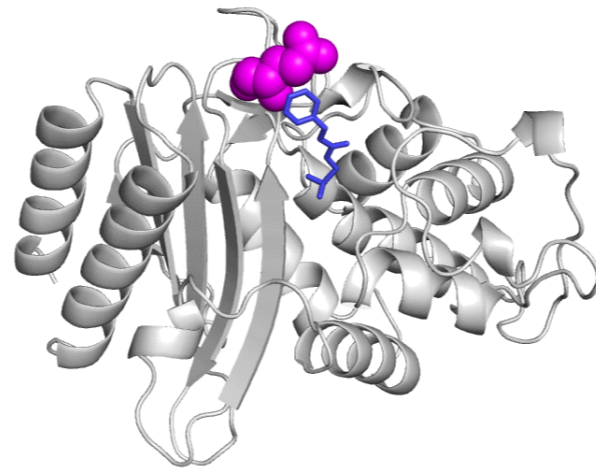
Structure

Function

Mutations impact protein function

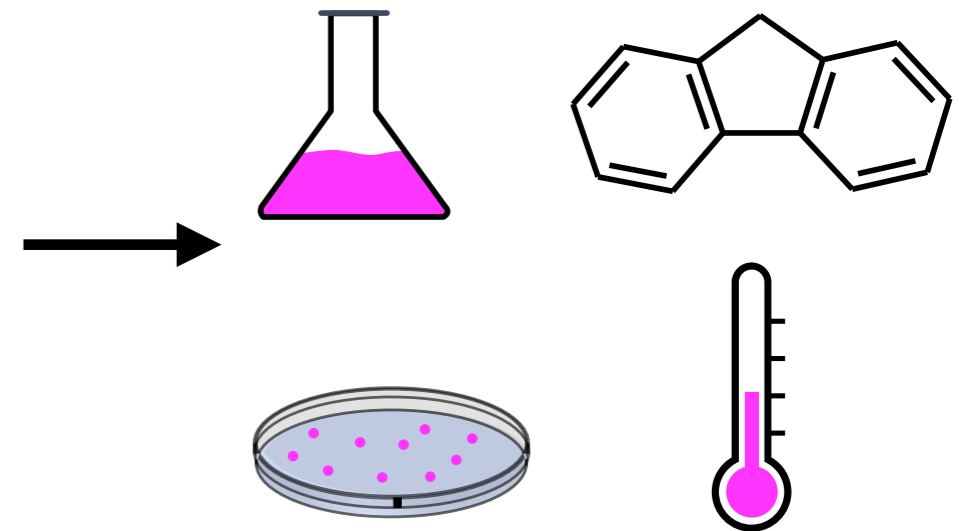
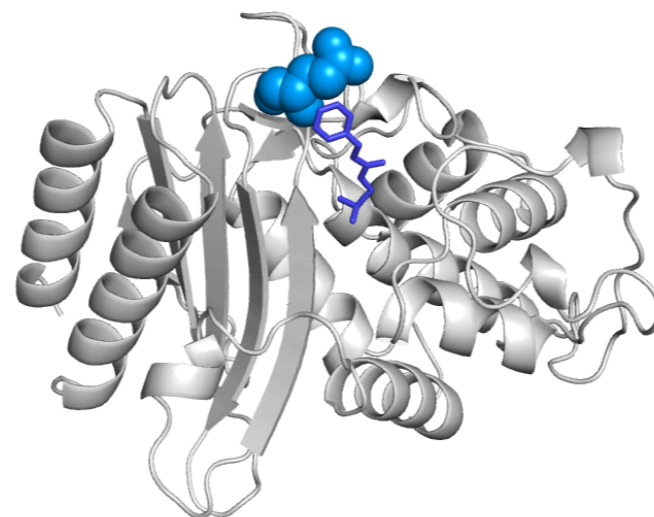
MSIQHFRVALIPFFAAFCCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAFLHNMGDHVTRLR
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



MSIQHFRVALIPFFAAFCCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAFLHNMGDHVTRLR
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → I



Sequence

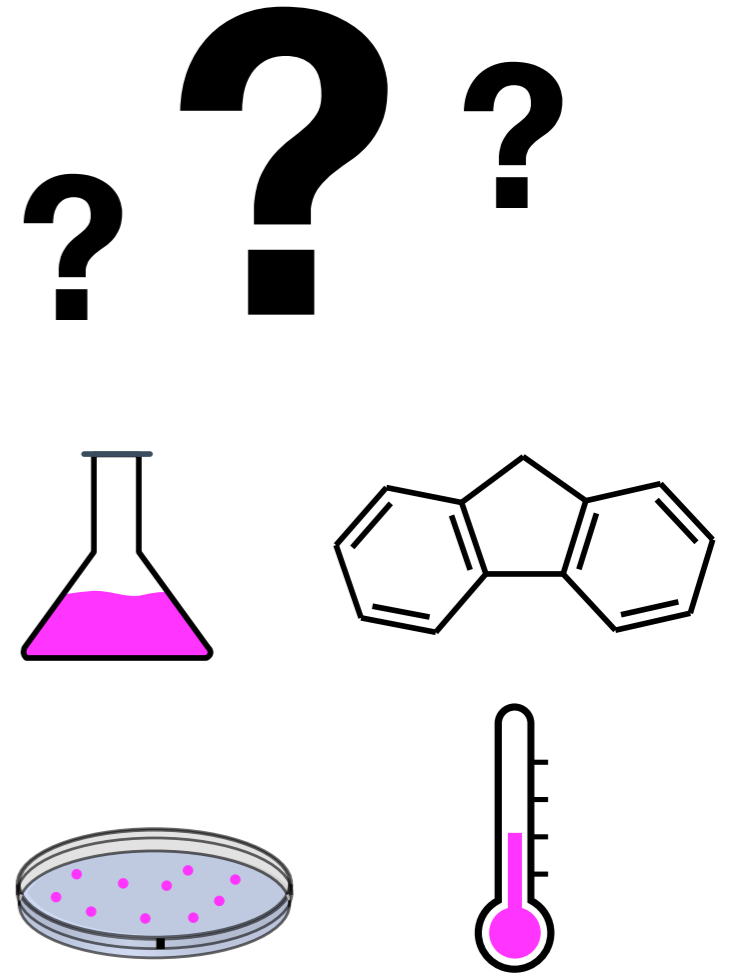
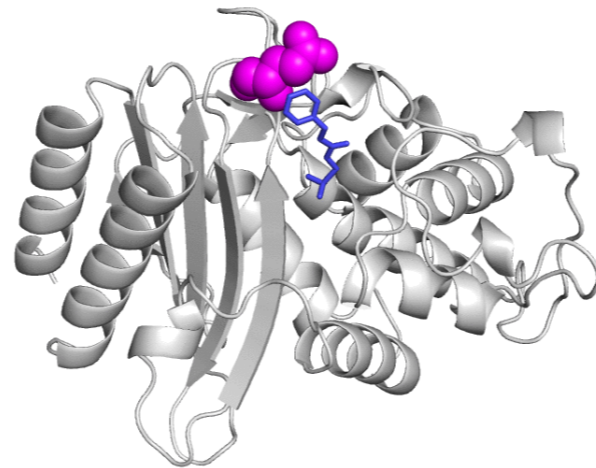
Structure

Function

Mutations impact protein function

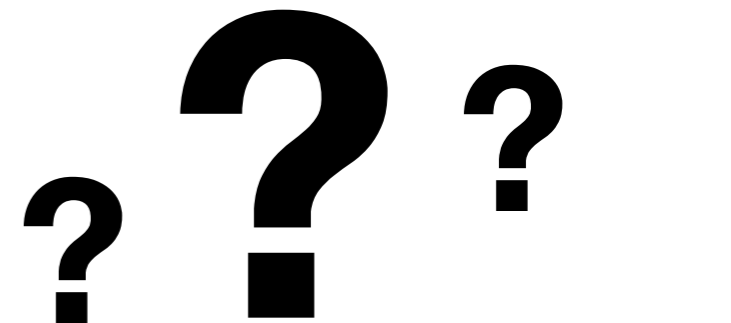
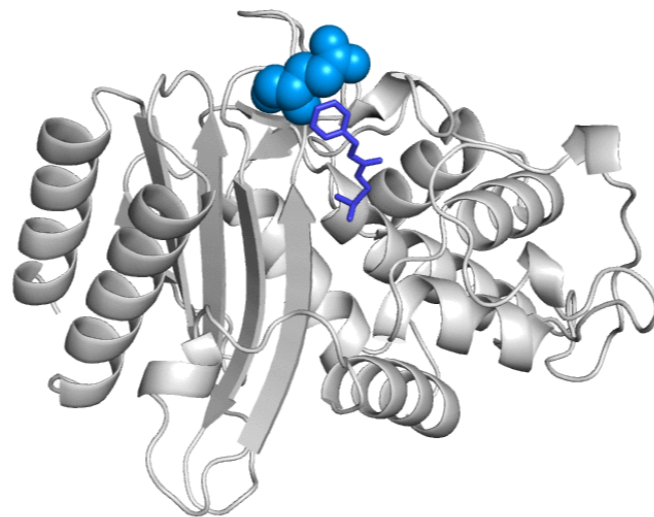
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLR
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L



MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLR
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGFLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → I



Sequence

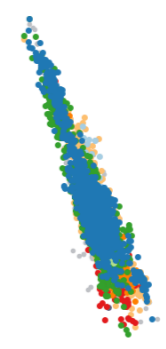
Structure

Function

High-Throughput Data Acquisition



Machine Learning

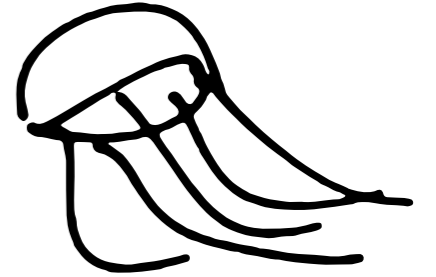
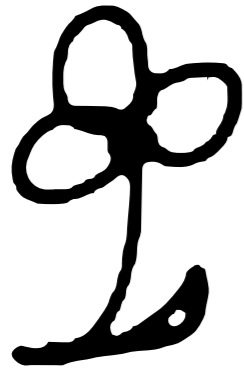


Update 10

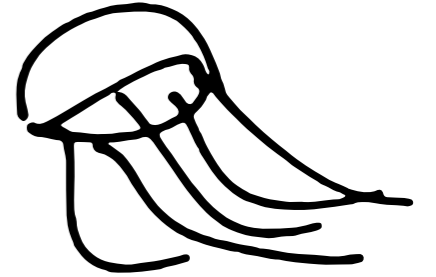
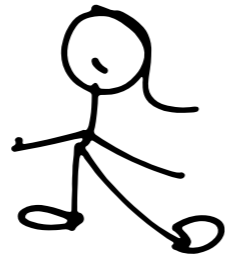
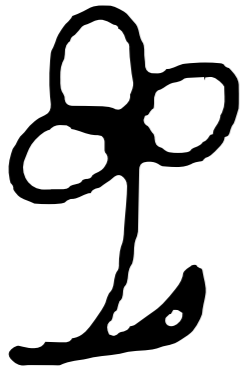
β -lactamase sequence family

- Acidobacteria
- Actinobacteria
- Bacteroidetes
- Chloroflexi
- Cyanobacteria
- Deinococcus-Thermus
- Firmicutes
- Fusobacteria
- Proteobacteria

DNA sequencing is becoming very cheap and easy

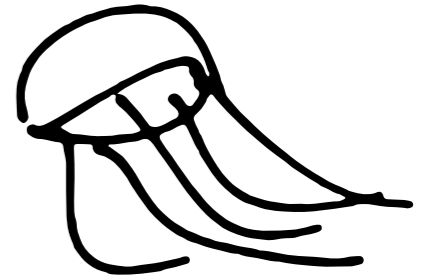
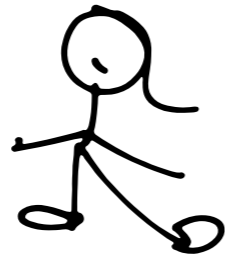
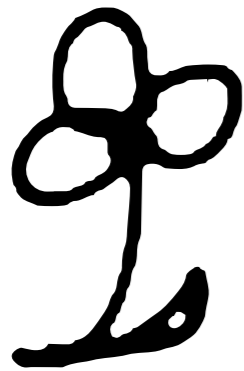


DNA sequencing is becoming very cheap and easy

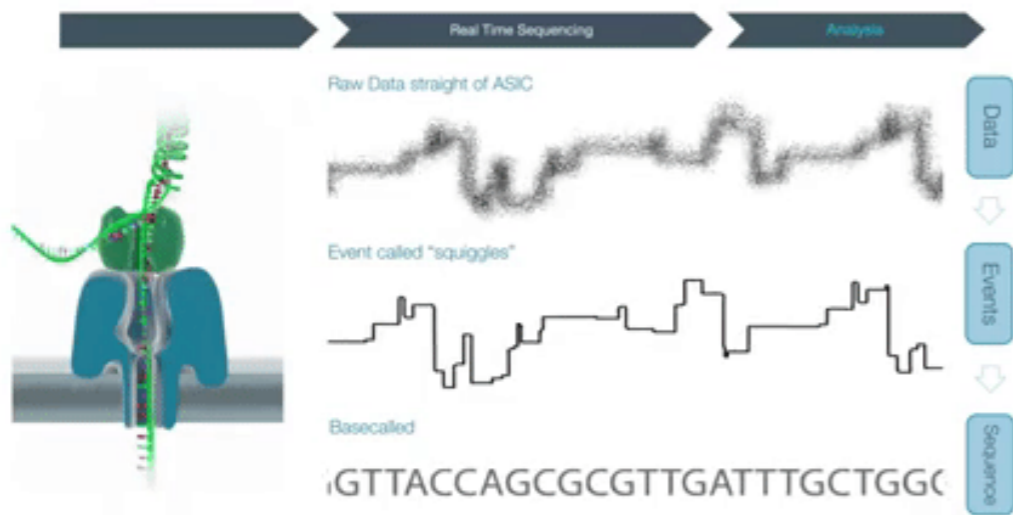


↓
Isolate
DNA

DNA sequencing is becoming very cheap and easy

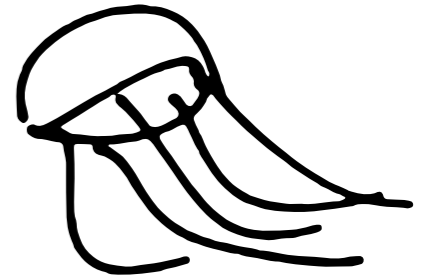
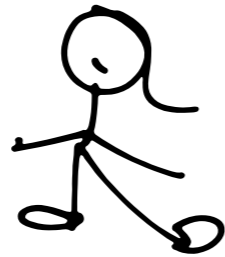
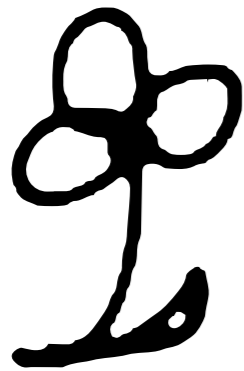


↓
Isolate
DNA



Sequence - Oxford Nanopore

DNA sequencing is becoming very cheap and easy



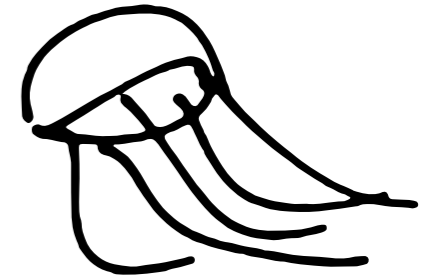
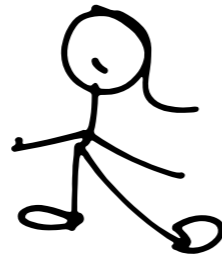
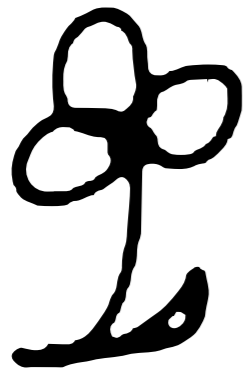
Isolate
DNA

Deposit in public
databases



Sequence - Oxford Nanopore

DNA sequencing is becoming very cheap and easy

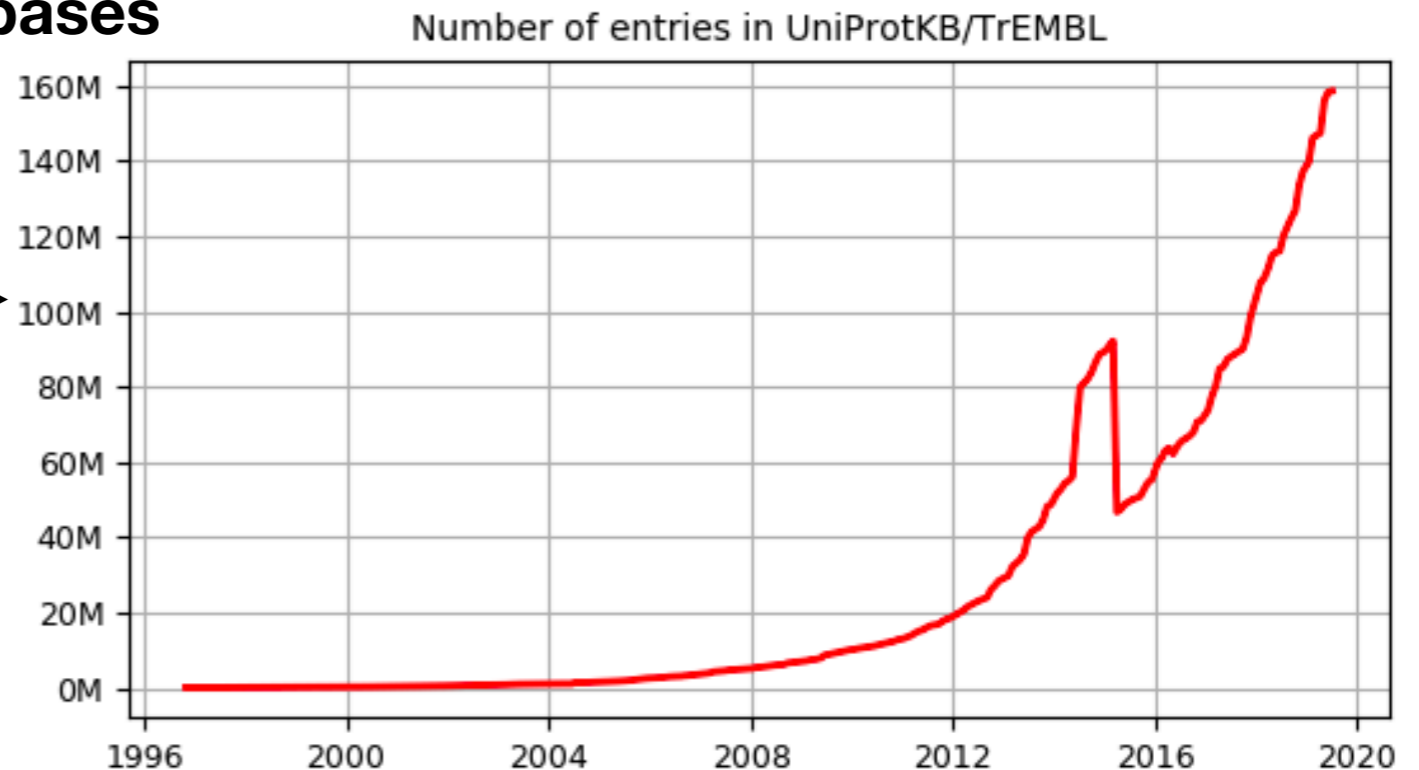


Isolate DNA

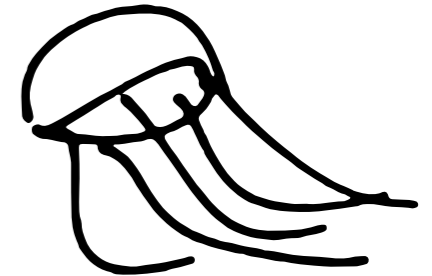
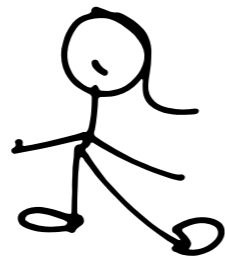
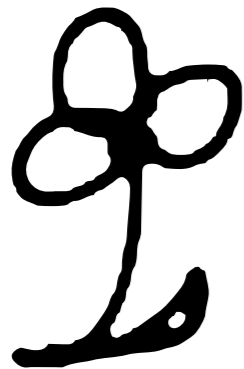
Deposit in public databases



Sequence - Oxford Nanopore

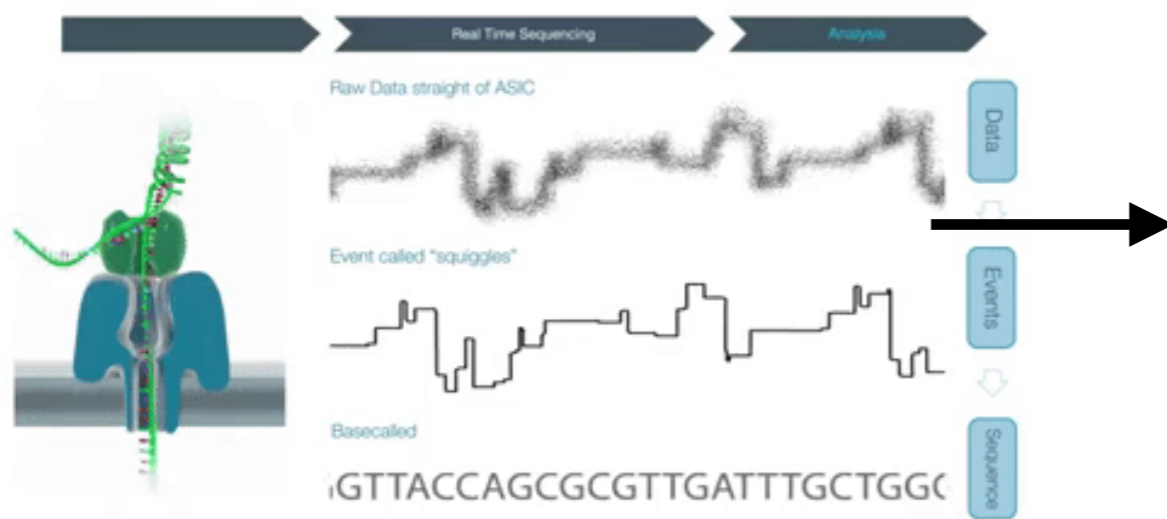


DNA sequencing is becoming very cheap and easy

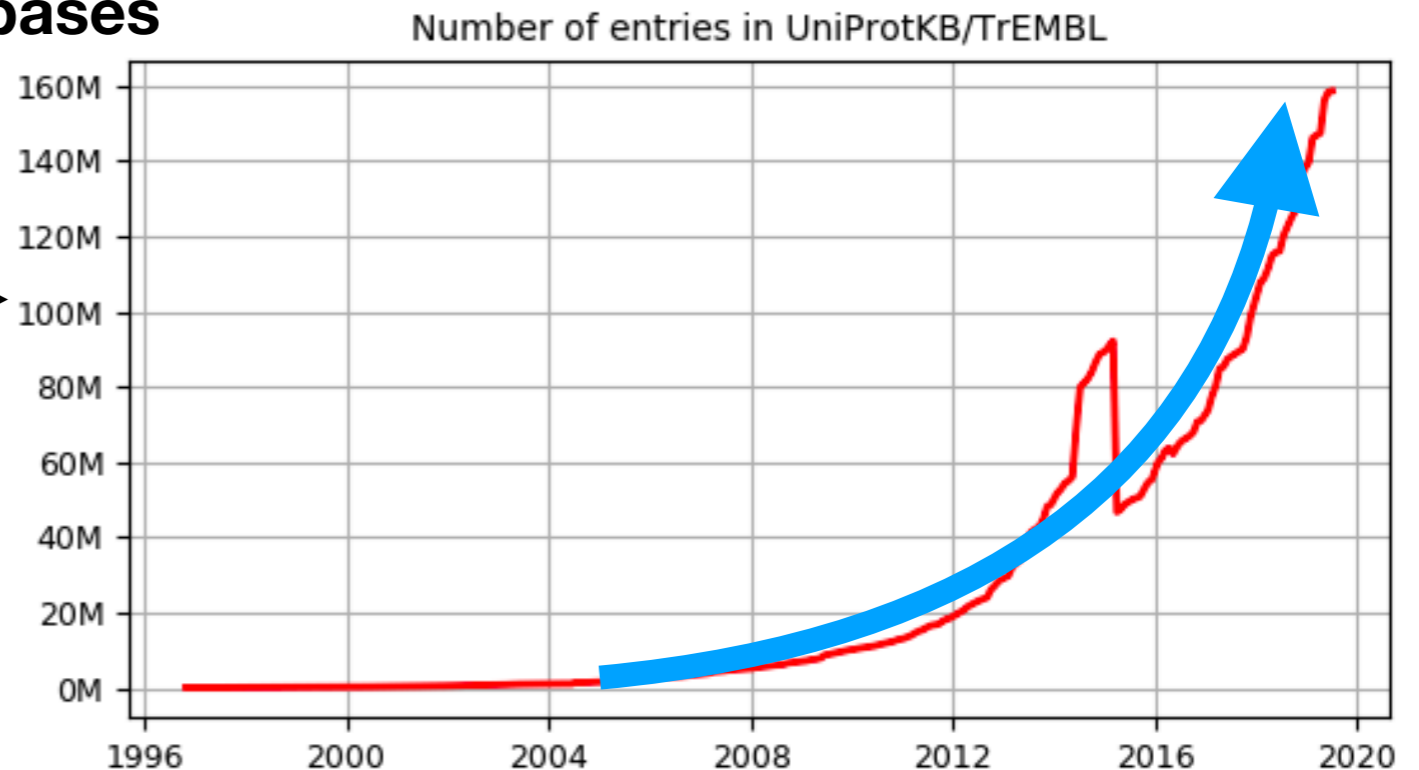


Isolate DNA

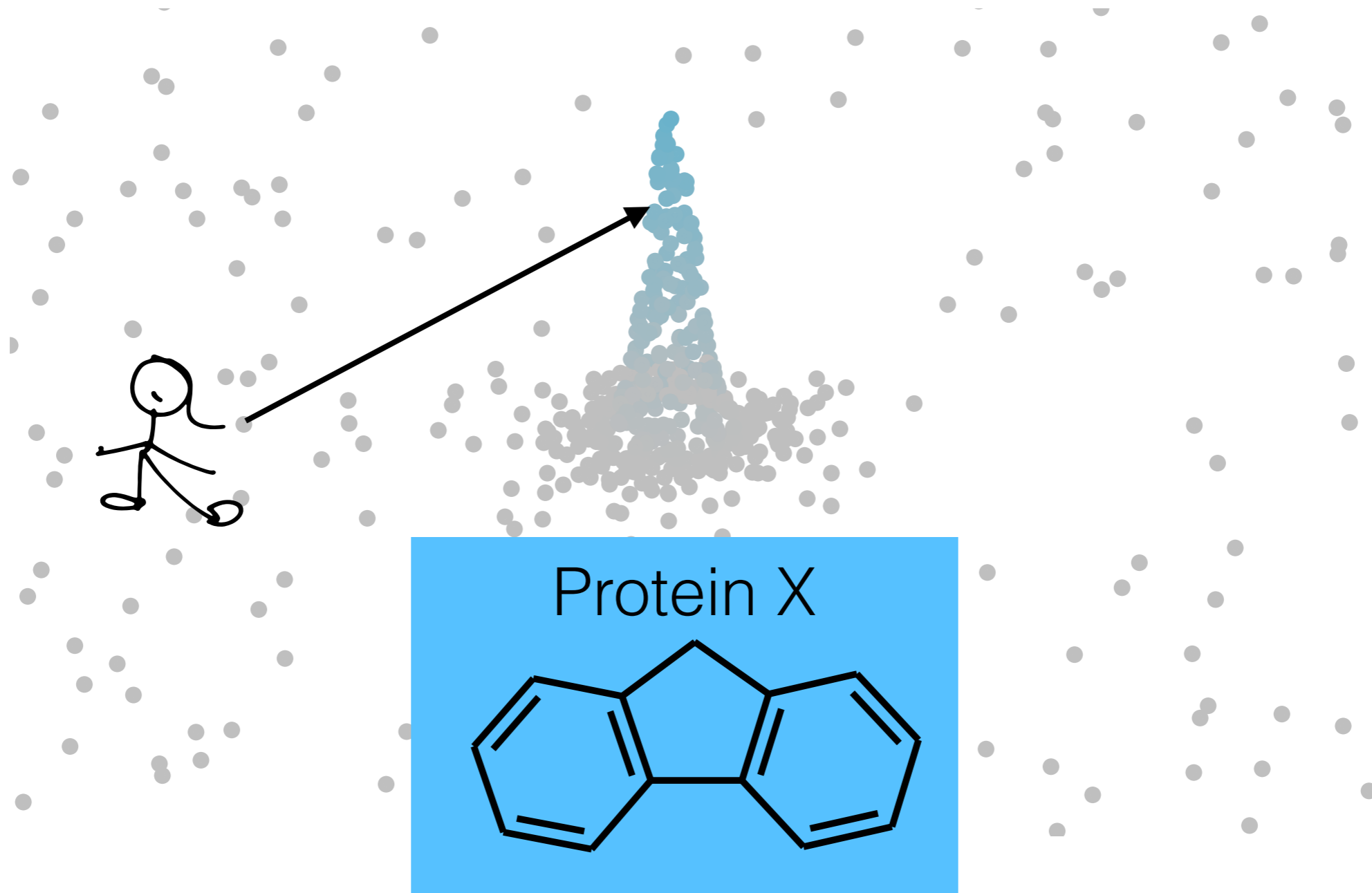
Deposit in public databases



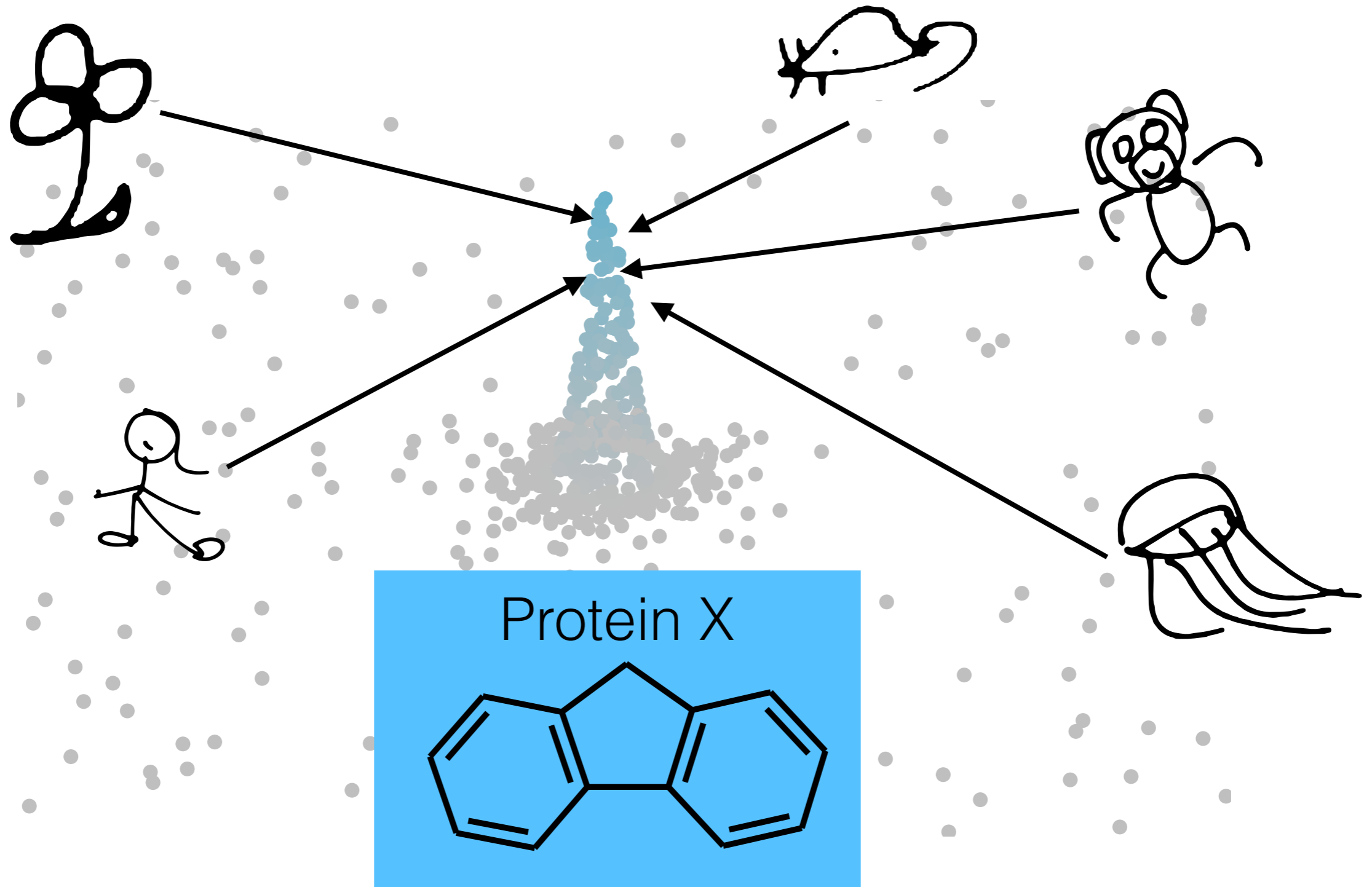
Sequence - Oxford Nanopore



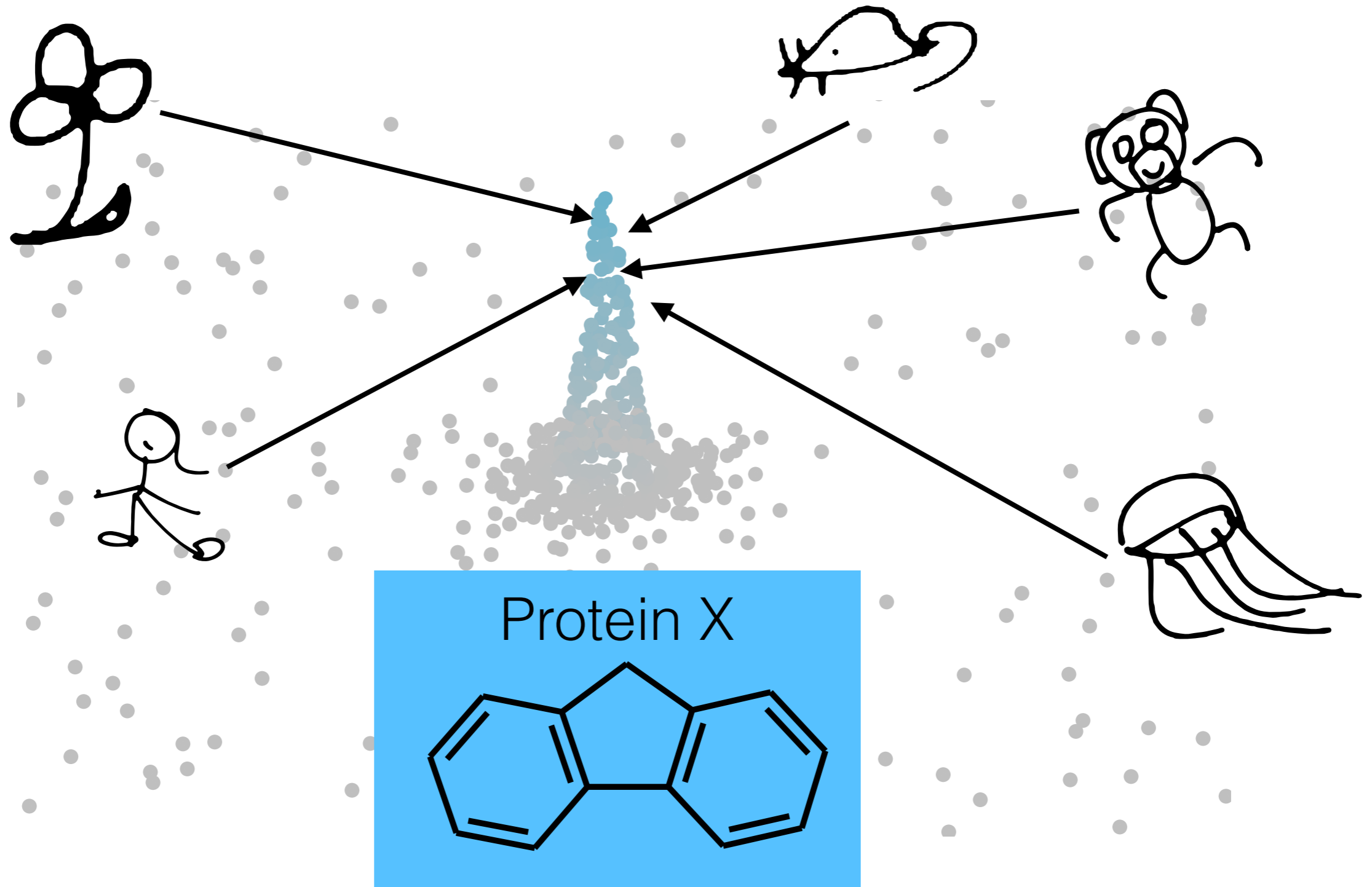
Lots of other examples of Protein X are available



Lots of other **examples** of **Protein X** are available



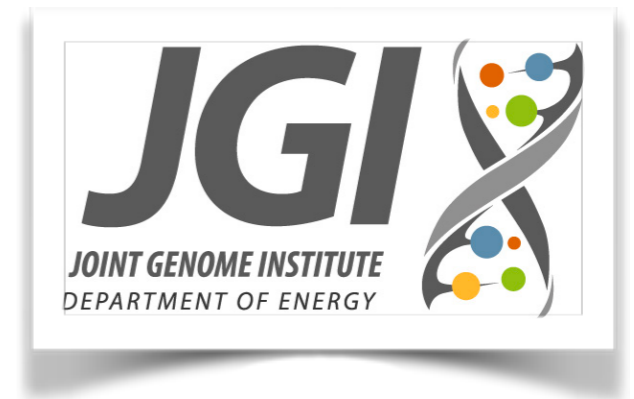
Lots of other **examples** of **Protein X** are **available**



All are **functional, homologous examples** of Protein X

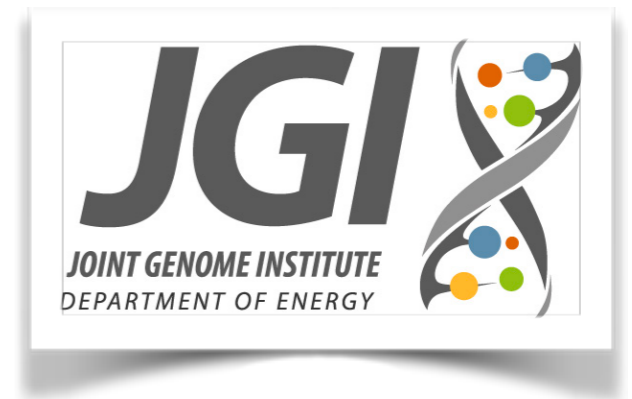
Lots of other **examples** of **Protein X** are **available**

Sequences are found in public
genome databases.



Lots of other **examples** of **Protein X** are **available**

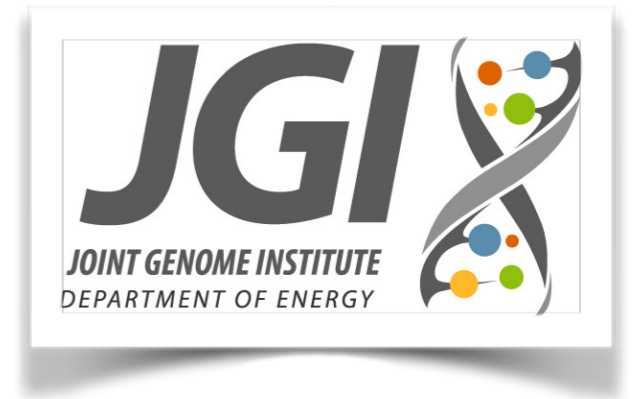
Sequences are found in public
genome databases.



Natural **evolution** is an **experiment**, in **parallel**.

Lots of other **examples** of **Protein X** are **available**

Sequences are found in public **genome databases.**



Natural **evolution** is an **experiment**, in **parallel**.

Assumption:

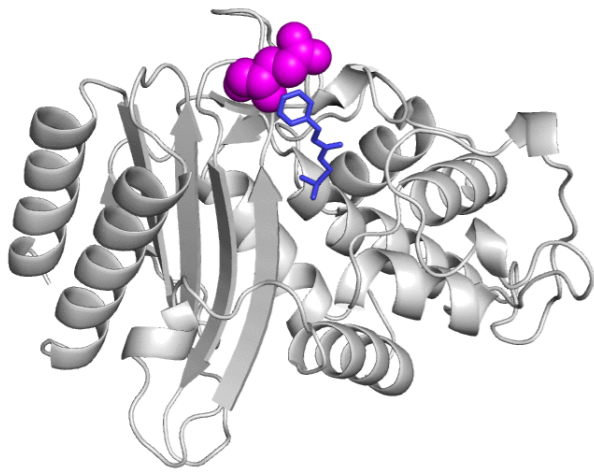
Present in database: Tolerated

Not in database: Deleterious

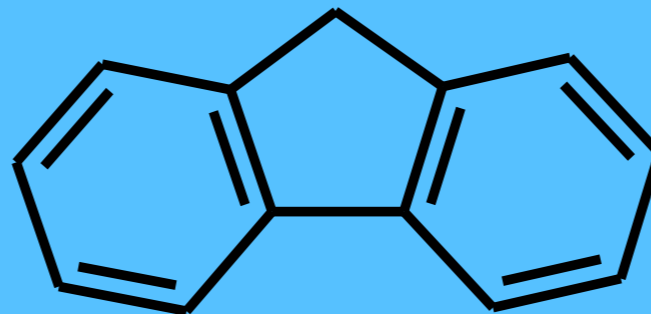
Can we predict which mutation will be tolerable?

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



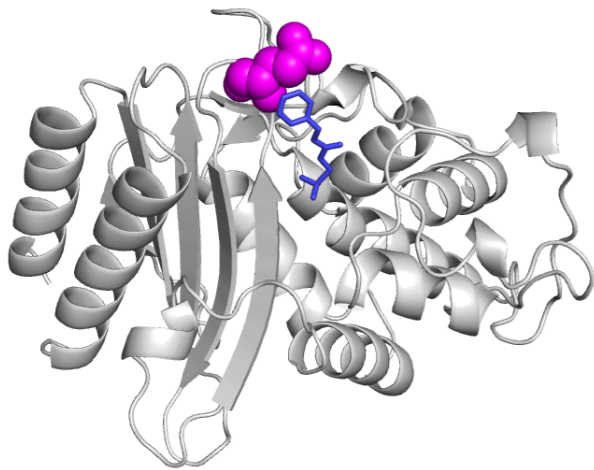
Protein X



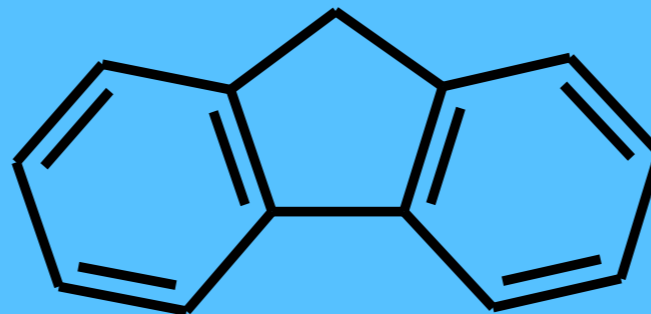
Can we **predict** which **mutation** will be **tolerable**?

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRLD
RWEPELNEAIPNDERDTTTPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



Protein X

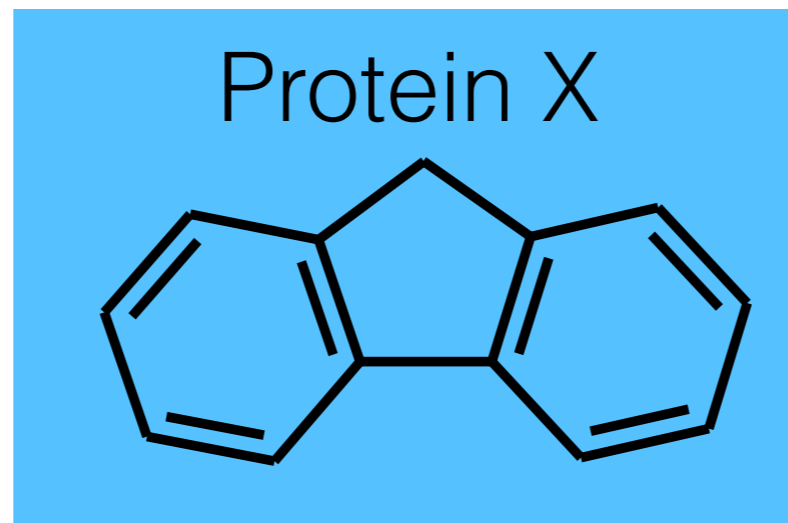
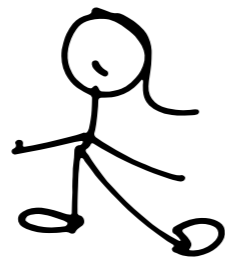
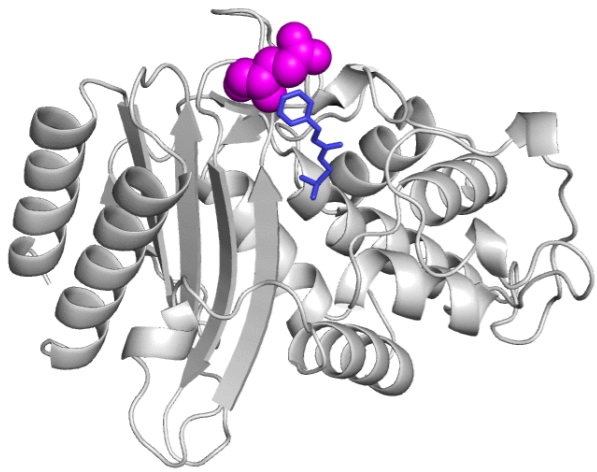


Can we predict which mutation will be tolerable?

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP L LRSALPAGWFIA
DKSGAG **E** RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L

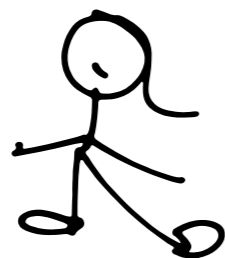
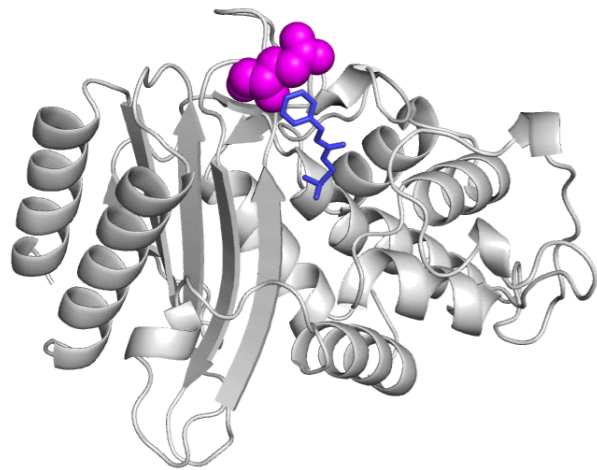
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP L LRSALPAGWFIA
DKSGAG **E** RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



Can we predict which mutation will be tolerable?

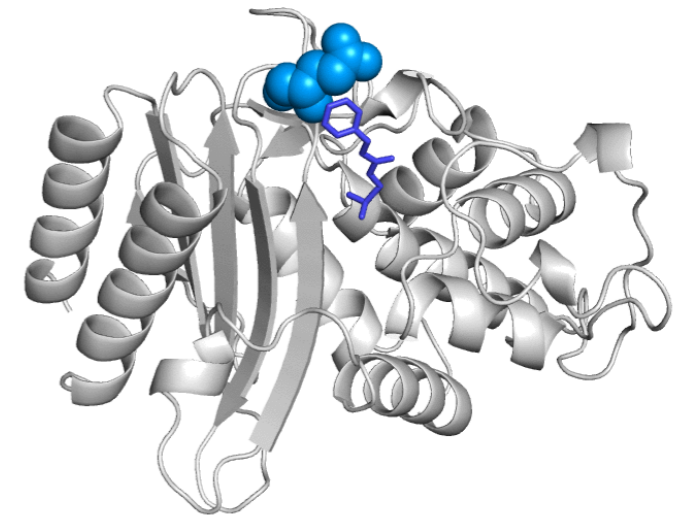
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKL LTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **E** RGS RGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L

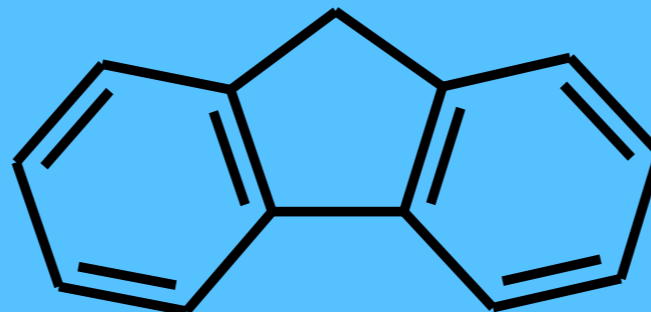


E → I

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKL LTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **E** RGS RGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW



Protein X



Can we predict which mutation will be tolerable?

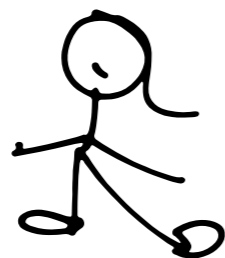
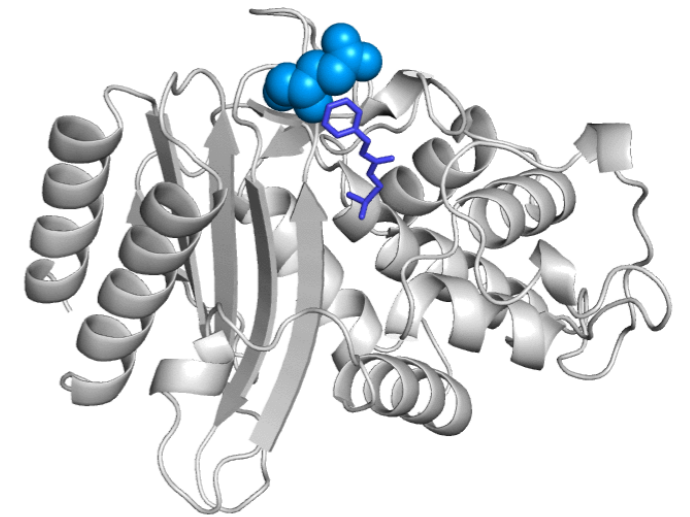
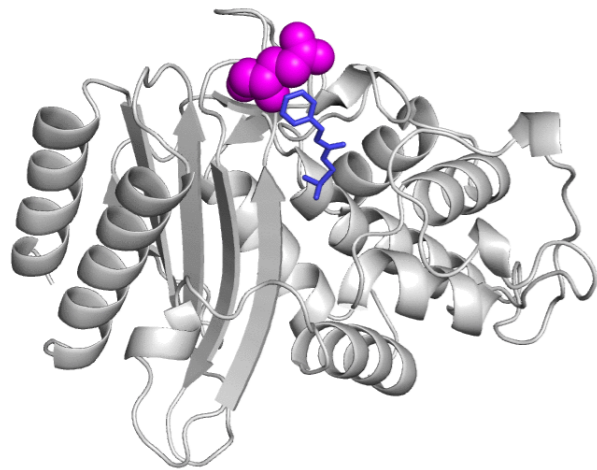
MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

E → L

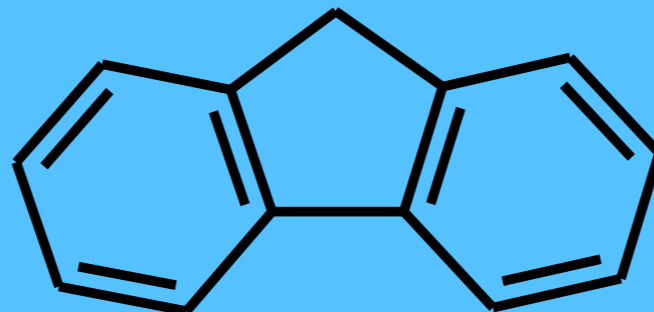
E → I

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGD HVTRL D
RWEPELNEAIPNDERDTTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGPLLRSALPAGWFIA
DKSGAG**E**RGSRGIIAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

???



Protein X



Can we predict which mutation will be tolerable?

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGDHSVTRLD
RWEPELNEAIPNDERD TTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **E** RGS RGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

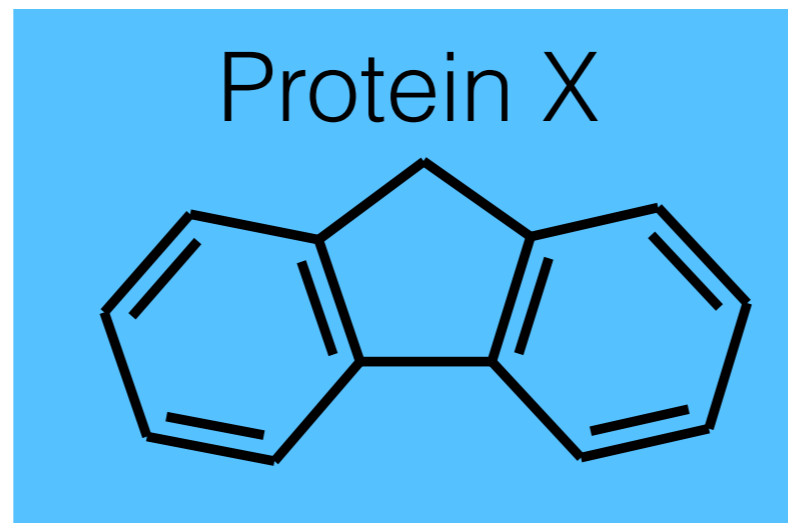
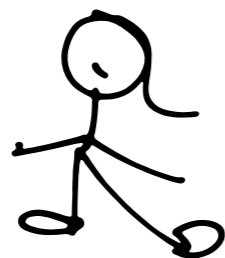
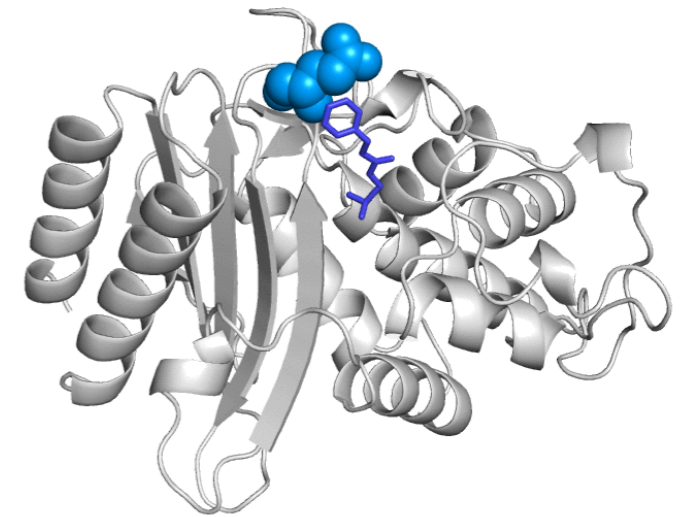
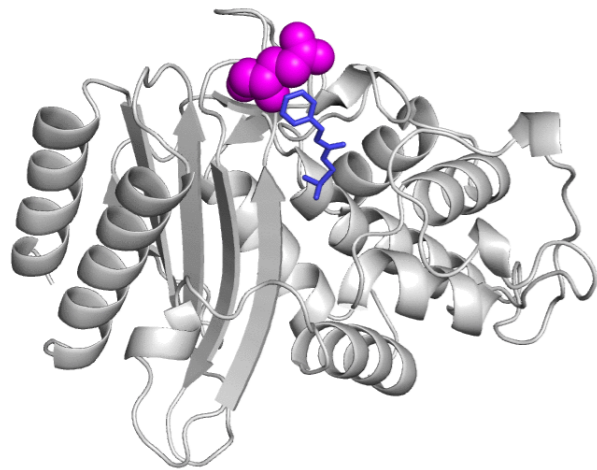
E → L

E → I

MSIQHFRVALIPFFAAFCLPVFA
HPETLVKVKDAEDQLGARVGYIE
LDLNSGKILESFRPEERFPMMST
FKVLLCGAVLSRVDAGQEQLGRR
IHYSQNDLVEYSPVTEKHLTDGM
TVRELCSAAITMSDNTAANLLLT
TIGGPKELTAF LHNMGDHSVTRLD
RWEPELNEAIPNDERD TTMPAAM
ATTLRKLLTGELLTLASRQQLID
WMEADKVAGP LLRSALPAGWFIA
DKSGAG **E** RGS RGI IAALGPDGKP
SRIVVIYTTGSQATMDERNRQIA
EIGASLIKHW

???

How can we formulate this problem?



A **generative model** finds **probable**
“words” based on **context**

A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the _____.

A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the _____.

banana



A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the _____.

banana



running



A **generative model** finds **probable** “words” based on **context**

English sentence

We are going to go to the _____.

banana



running




zoo



A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the zoo.



P(x) finds
probable “words”
based on **context**

A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the zoo.

A diagram illustrating context windows for the sentence "We are going to go to the zoo.". The sentence is written in black text. Below the text, several black curved lines represent context windows. One window is under "We", another under "are", and a larger one under "going to go to the zoo.". A single window is also shown above "go to the zoo.".

P(x) finds
probable “words”
based on **context**

Protein sequence

AQKLYLTHIDAEV**EGADTL**FITEVKQVF

A diagram illustrating context windows for the protein sequence "AQKLYLTHIDAEV**EGADTL**FITEVKQVF". The sequence is written in black text, with "EGADTL" highlighted in blue. Below the text, several black curved lines represent context windows. One window is under "AQKLYL", another under "THIDAEV", and a larger one under "EGADTL FITEVKQVF". A single window is also shown above "EGADTL FITEVKQVF".

A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the zoo.



Protein sequence

AQKLYLTHIDAEV**EGADTL**FITEVKQVF



A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the zoo.



The diagram shows arcs connecting words in the sentence: 'We' to 'are', 'going' to 'to', 'to' to 'go', 'go' to 'to', 'to' to 'the', and 'the' to 'zoo'.

Protein sequence

AQKLYLTHIDAEV**E**GADTLFITEV**K**QVF




The diagram shows arcs connecting amino acids in the sequence: 'A' to 'K', 'L' to 'Y', 'L' to 'T', 'T' to 'H', 'H' to 'I', 'I' to 'D', 'D' to 'A', 'A' to 'E', 'E' to 'V', 'V' to 'G', 'G' to 'A', 'A' to 'D', 'D' to 'T', 'T' to 'L', 'L' to 'F', 'F' to 'I', 'I' to 'T', 'T' to 'E', 'E' to 'V', 'V' to 'K', and 'K' to 'Q'. A red 'H' is positioned above the gap after 'QVF', with a line connecting it to the end of the sequence.



A **generative model** finds **probable**
“words” based on **context**

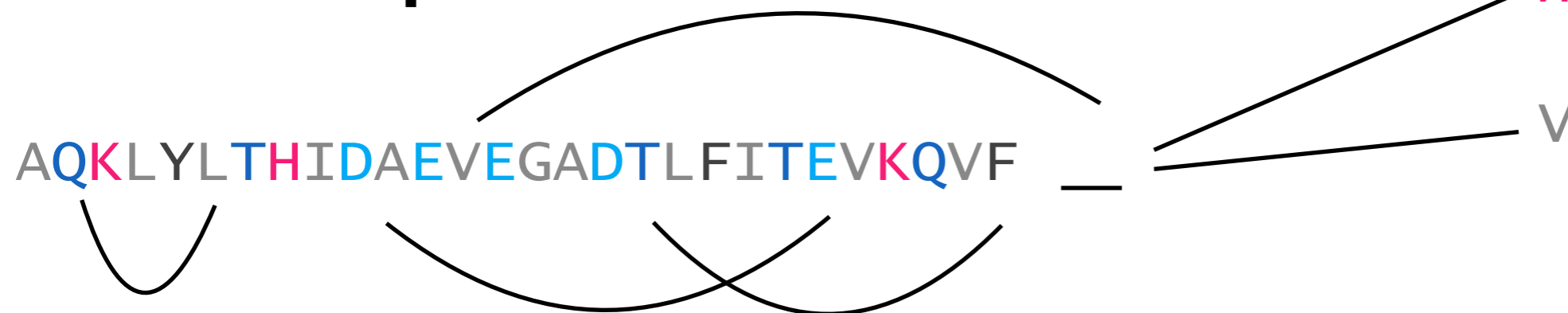
English sentence

We are going to go to the zoo.



Protein sequence

AQKLYLTHIDAEV**EGADTL**FITEVKQVF



A **generative model** finds **probable**
“words” based on **context**

English sentence

We are going to go to the zoo.

Protein sequence

AQKLYLTHIDAEV**EGADTL**FITEVKQVF

H



V



T



Neural networks power generative models

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

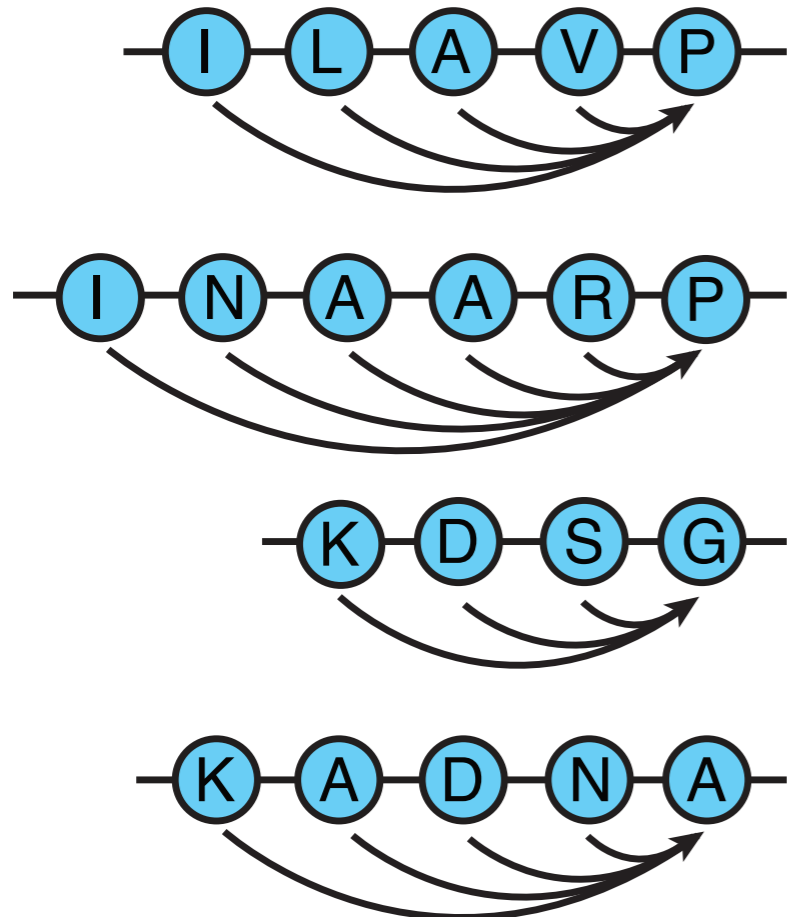
We are going to go to the zoo.



Neural networks power generative models

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

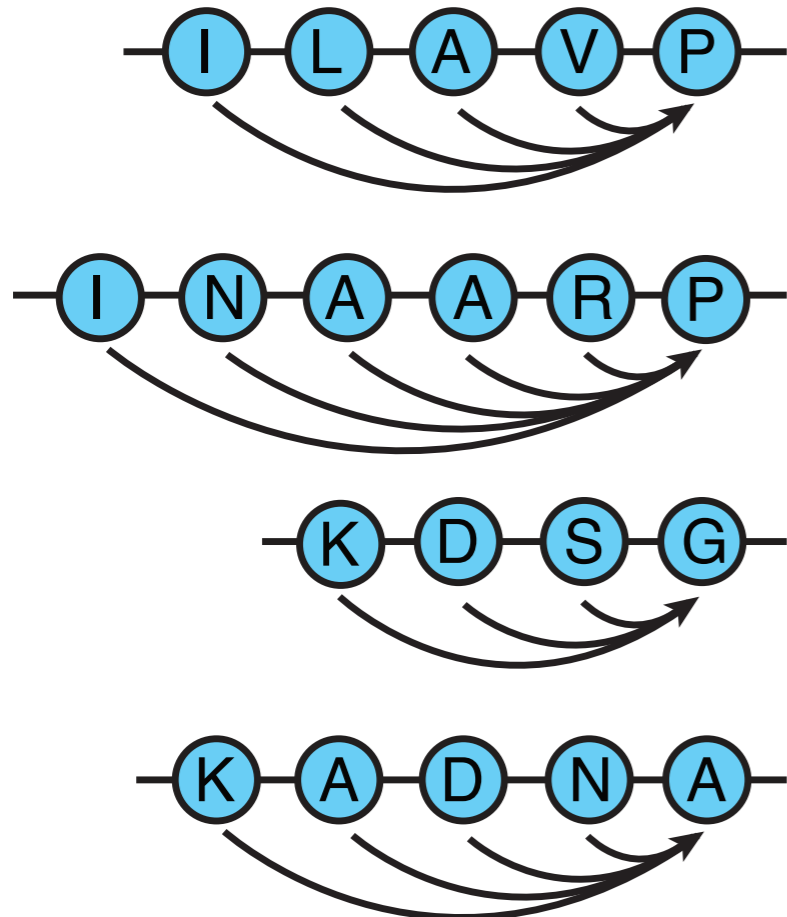
Autoregressive model



Neural networks power generative models

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

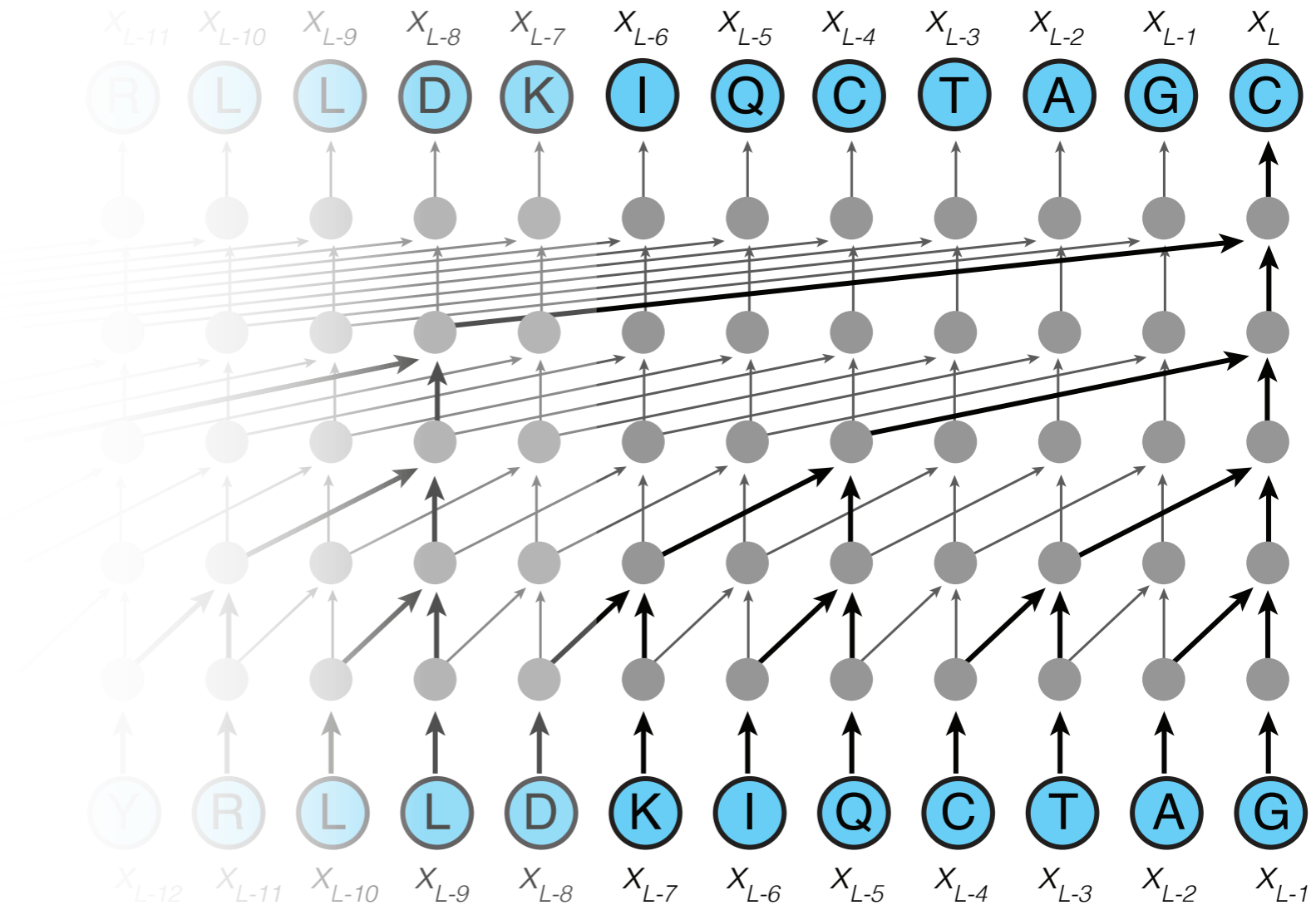
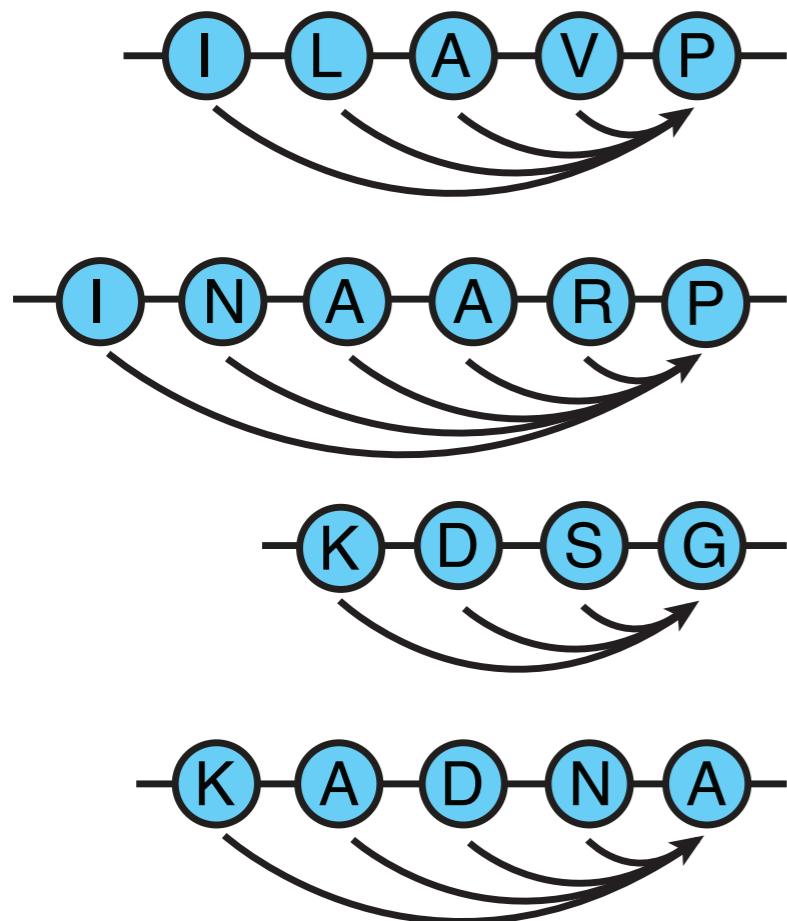
Autoregressive model



Neural networks power generative models

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

Autoregressive model

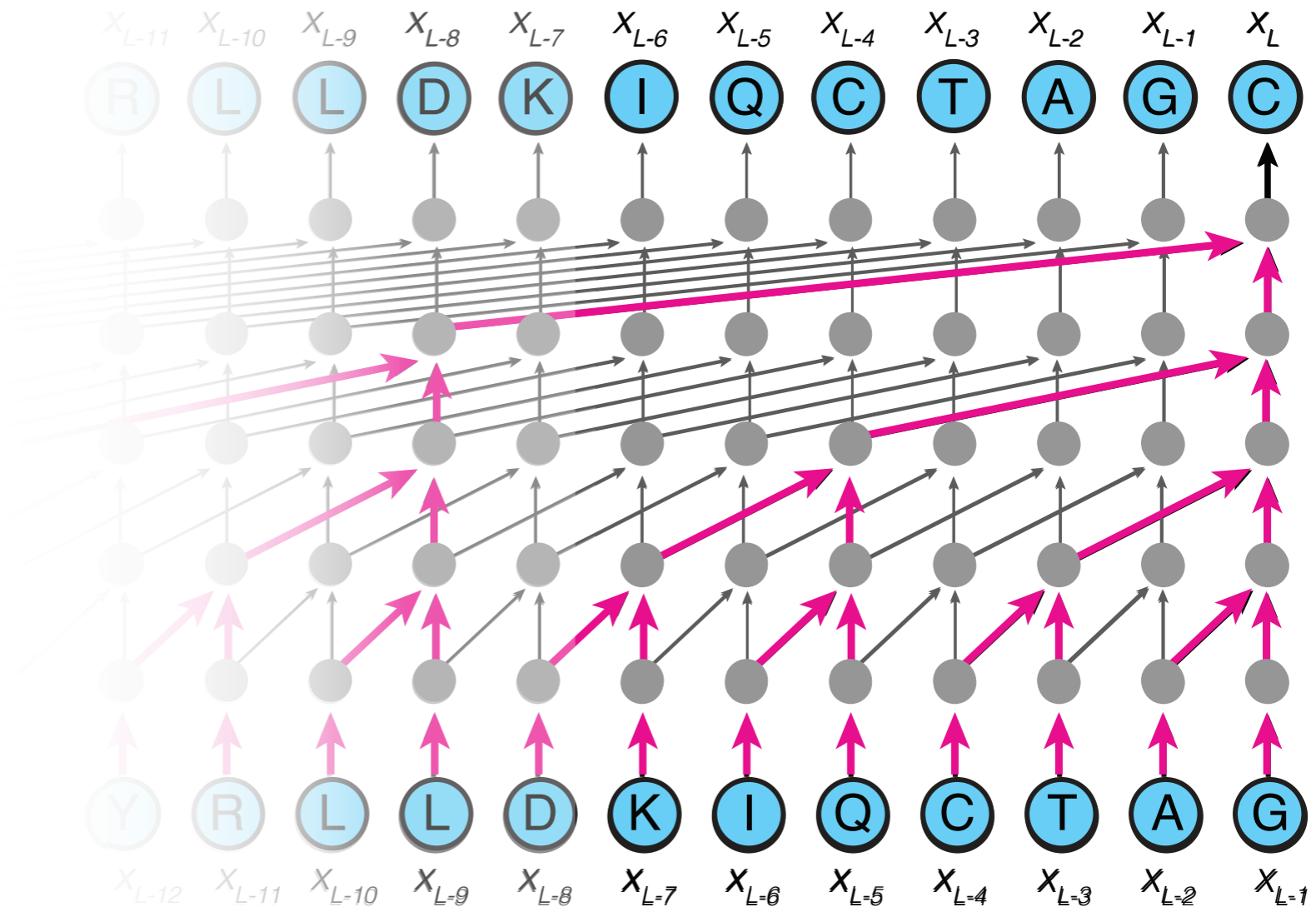
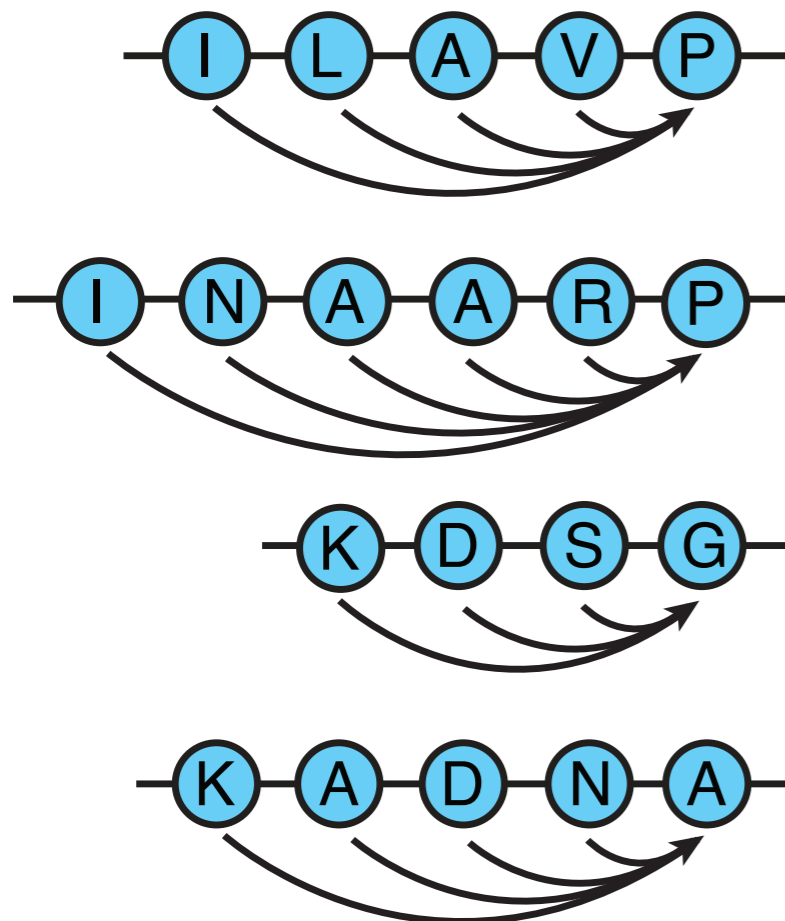


Dilated convolutional neural network

Neural networks power generative models

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

Autoregressive model



Dilated convolutional neural network

Utilizing an autoregressive likelihood

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$

Output **T** A **D** **T** L F I **T**

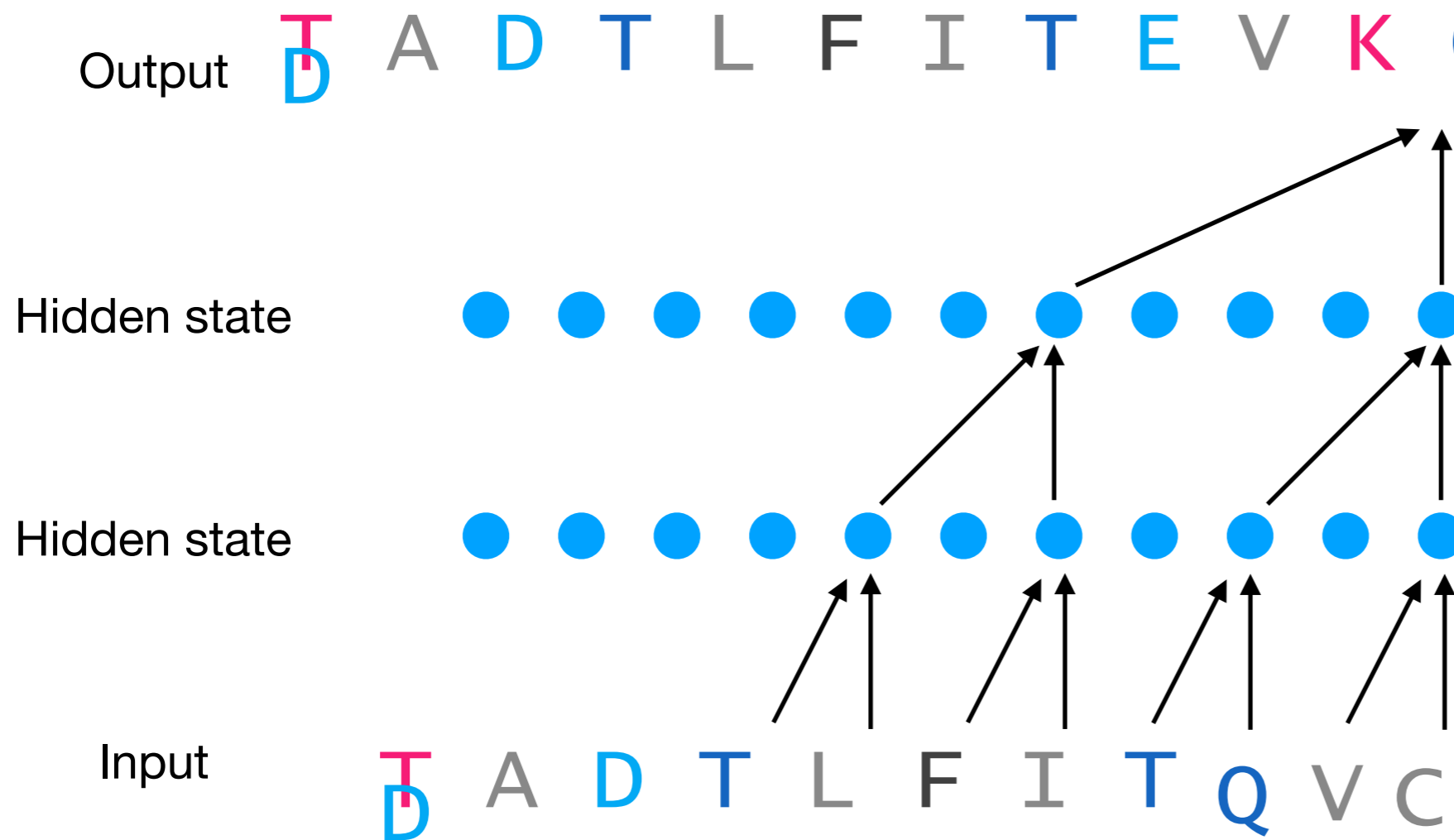
Hidden state ● ● ● ● ● ● ●

Hidden state ● ● ● ● ● ● ●

Input **T** A **D** **T** L F I **T**

Utilizing an autoregressive likelihood

$$p(x|\theta) = \prod_i^L p(x_i|x_{<i}, \theta)$$



Mutation effect prediction with an unsupervised model

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFKGD



$$p(\mathbf{x}|\theta)$$

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio**
for each mutant

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFECD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFECD
ADRLYMTKIHHEFECD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFECD
ADRLYLTQIRNKFKGD



$p(\mathbf{x}|\theta)$

$$\log \frac{p(\mathbf{x}_{\text{mut}}|\theta)}{p(\mathbf{x}_{\text{wild}}|\theta)}$$

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFECD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFECD
ADRLYMTKIHHEFECD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFECD
ADRLYLTQIRNKFKGD

↓
 $p(\mathbf{x}|\theta)$

2) Compute **Log Ratio**
for each mutant

$$\log \frac{p(\mathbf{x}_{\text{mut}}|\theta)}{p(\mathbf{x}_{\text{wild}}|\theta)}$$

“How much does this mutation look like what we’ve seen in nature?”

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio** for each mutant

Uses public data (effectively free)

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD

$p(\mathbf{x}|\theta)$

$$\log \frac{p(\mathbf{x}_{\text{mut}}|\theta)}{p(\mathbf{x}_{\text{wild}}|\theta)}$$

“How much does this mutation look like what we’ve seen in nature?”

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio** for each mutant

Uses public data (effectively free)

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD

Fast

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \theta)}{p(\mathbf{x}_{\text{wild}} | \theta)}$$

↓

$$p(\mathbf{x} | \theta)$$

“How much does this mutation look like what we’ve seen in nature?”

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio** for each mutant

Uses public data (effectively free)

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD

Fast

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \theta)}{p(\mathbf{x}_{\text{wild}} | \theta)}$$

Works on almost any protein

↓
 $p(\mathbf{x} | \theta)$

“How much does this mutation look like what we’ve seen in nature?”

Mutation effect prediction with an unsupervised model

1) Infer a **generative model** of the family

2) Compute **Log Ratio** for each mutant

Uses public data (effectively free)

AQKLYLTHIDAEVEGD
ADRLYMTKIHQFDGD
ADTLFITEVKQVFEGD
ADRLYMTKIHHTFDGD
ADKLYCTLIHNSFEGD
ADRLYMTKIHHEFEGD
ADRLYLTMIHQKFEAD
TDRLYITHIDETFEGD
ADRLYLTQIRNKFEGD

Fast

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \theta)}{p(\mathbf{x}_{\text{wild}} | \theta)}$$

Works on almost any protein

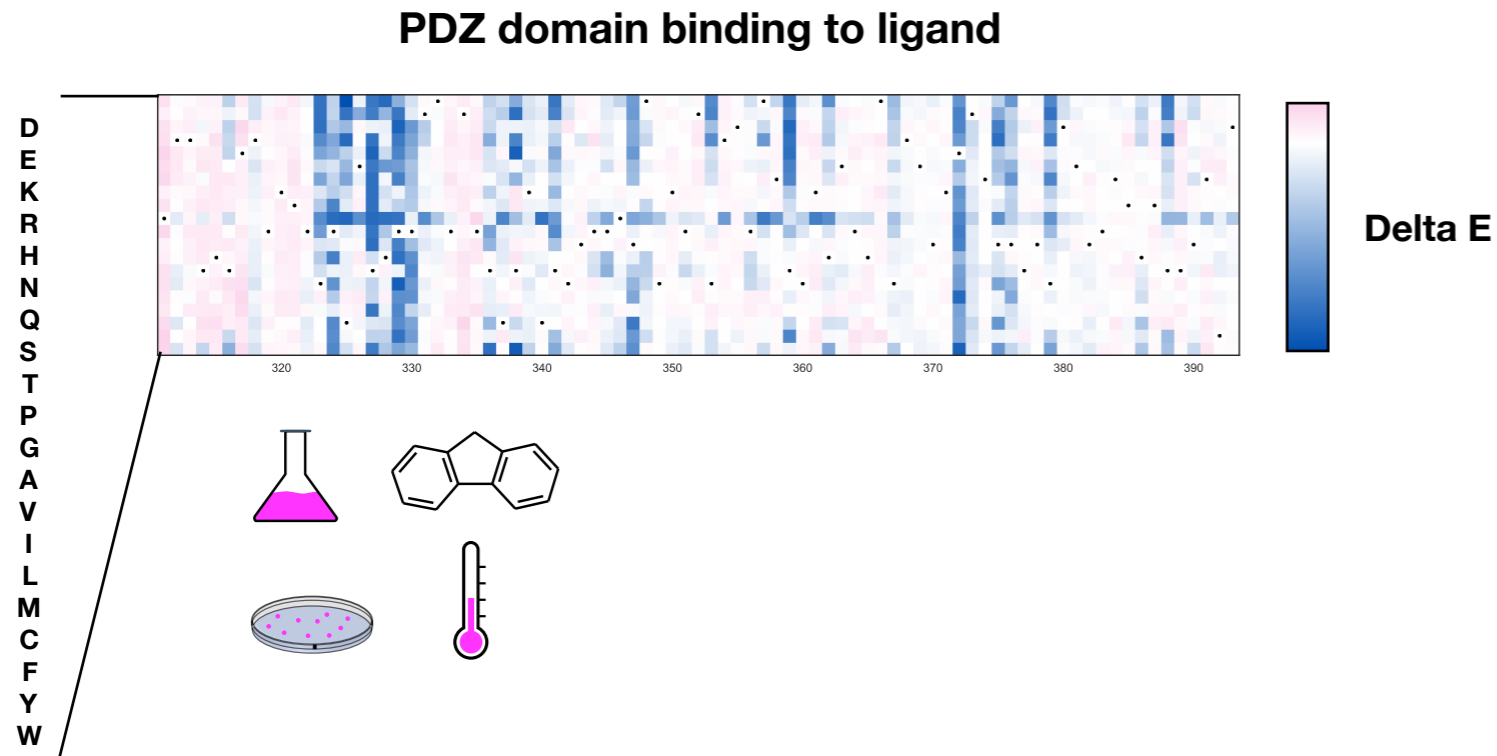
↓
 $p(\mathbf{x} | \theta)$

Accurate

“How much does this mutation look like what we’ve seen in nature?”

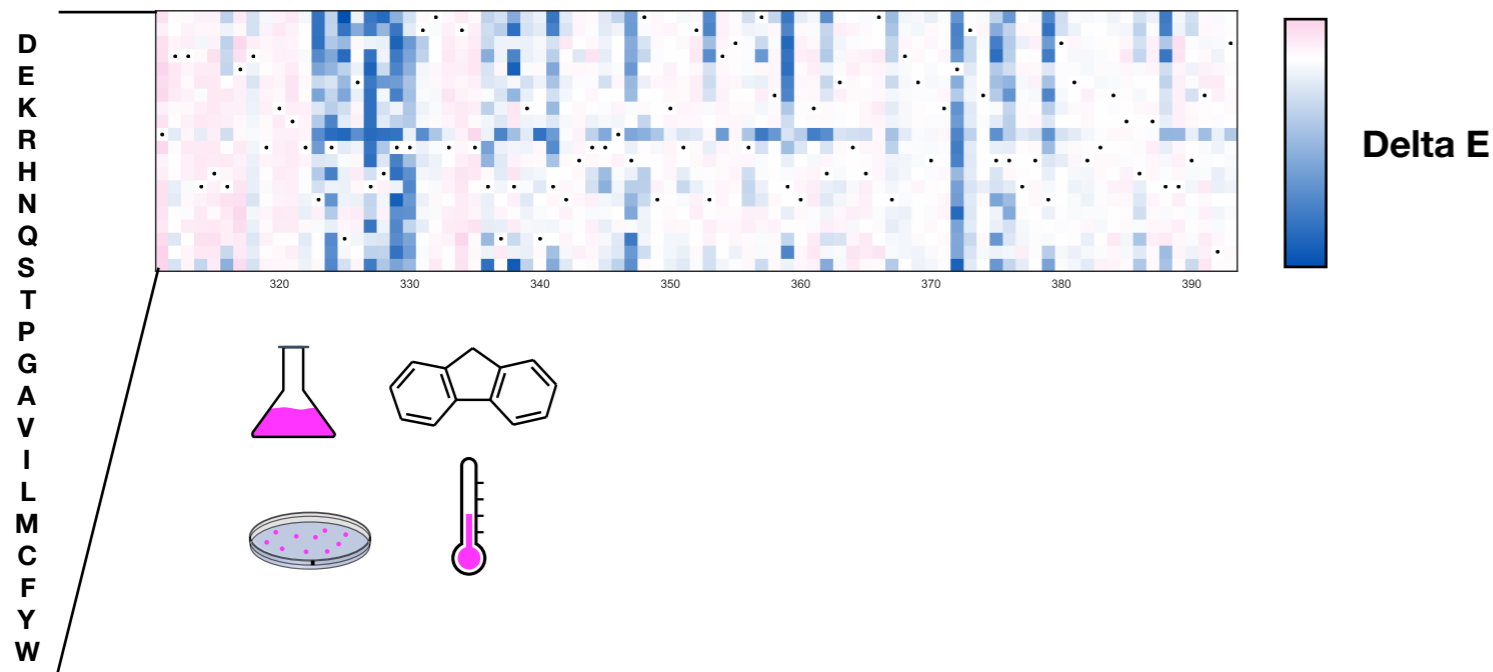
To evaluate performance, we collected ~30
saturation mutagenesis experiments

To evaluate performance, we collected ~30 **saturation mutagenesis experiments**

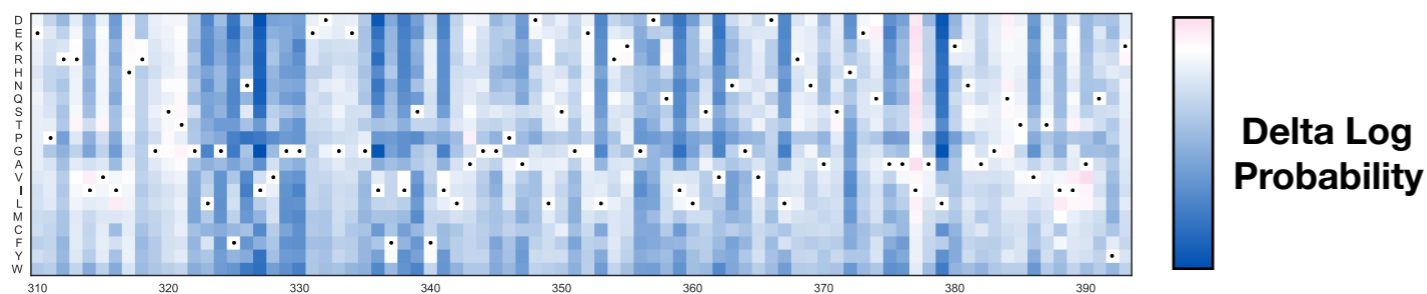


To evaluate performance, we collected ~30 **saturation mutagenesis experiments**

PDZ domain binding to ligand

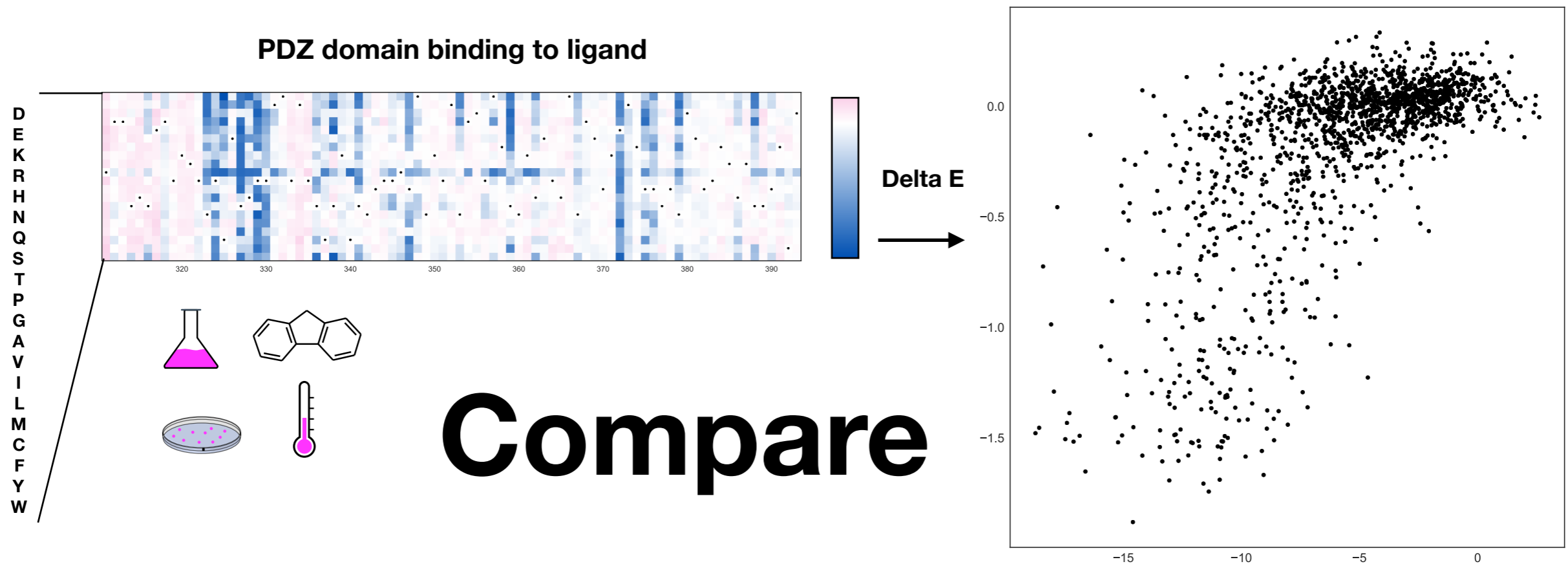


Log probability ratio



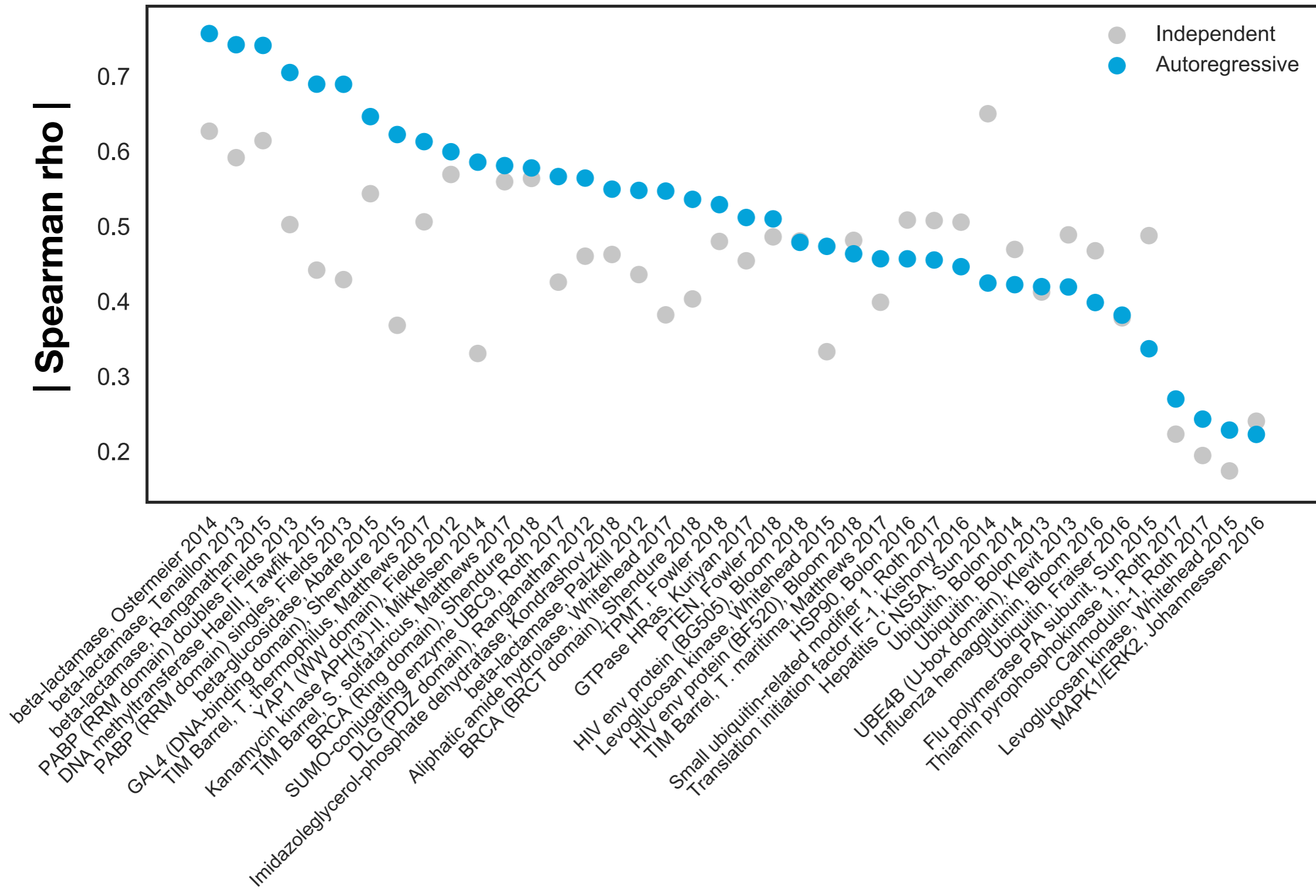
$$\log \frac{p(\mathbf{x}_{\text{mut}}|\boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}}|\boldsymbol{\theta})}$$

To evaluate performance, we collected ~30 saturation mutagenesis experiments

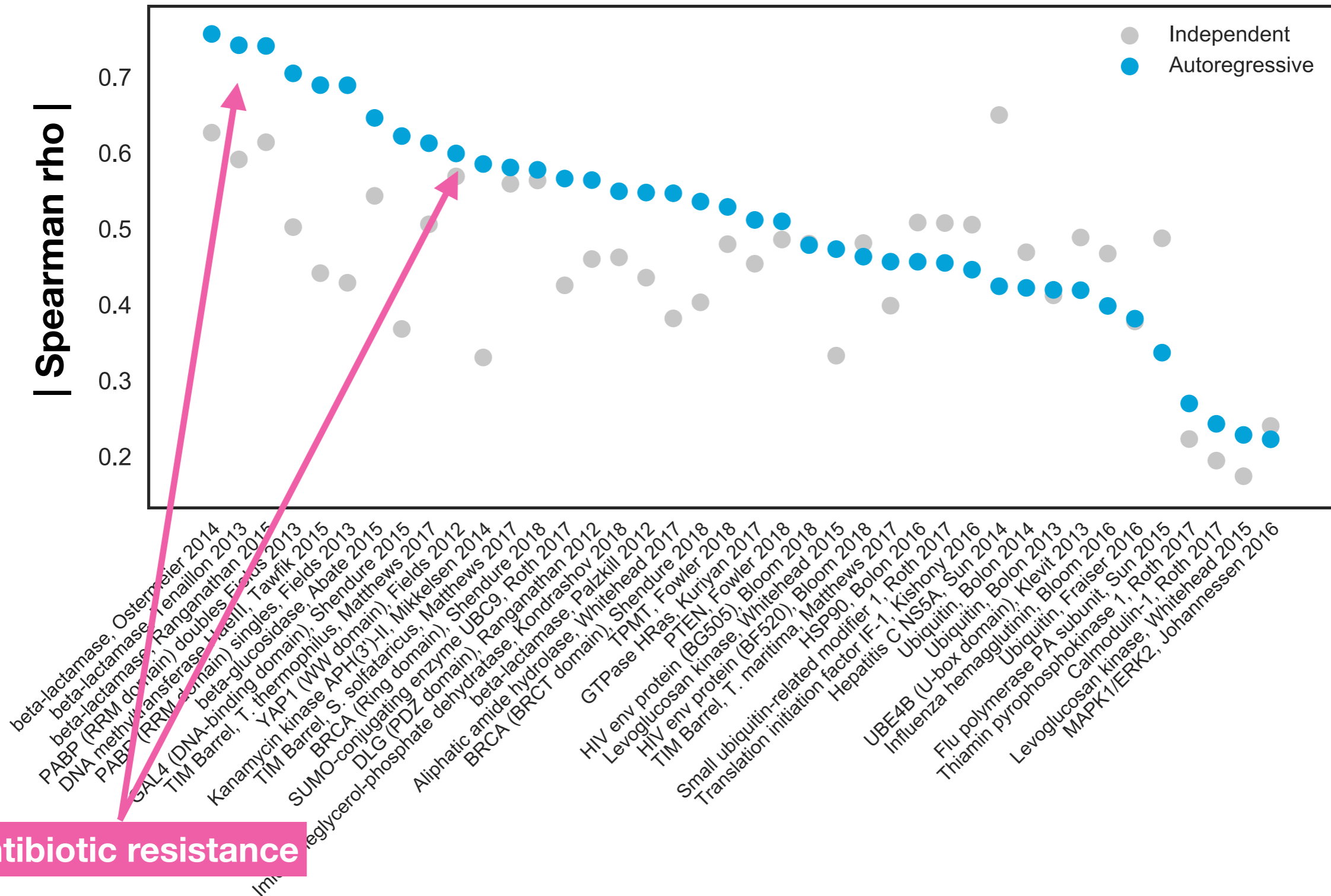


$$\log \frac{p(\mathbf{x}_{\text{mut}}|\boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}}|\boldsymbol{\theta})}$$

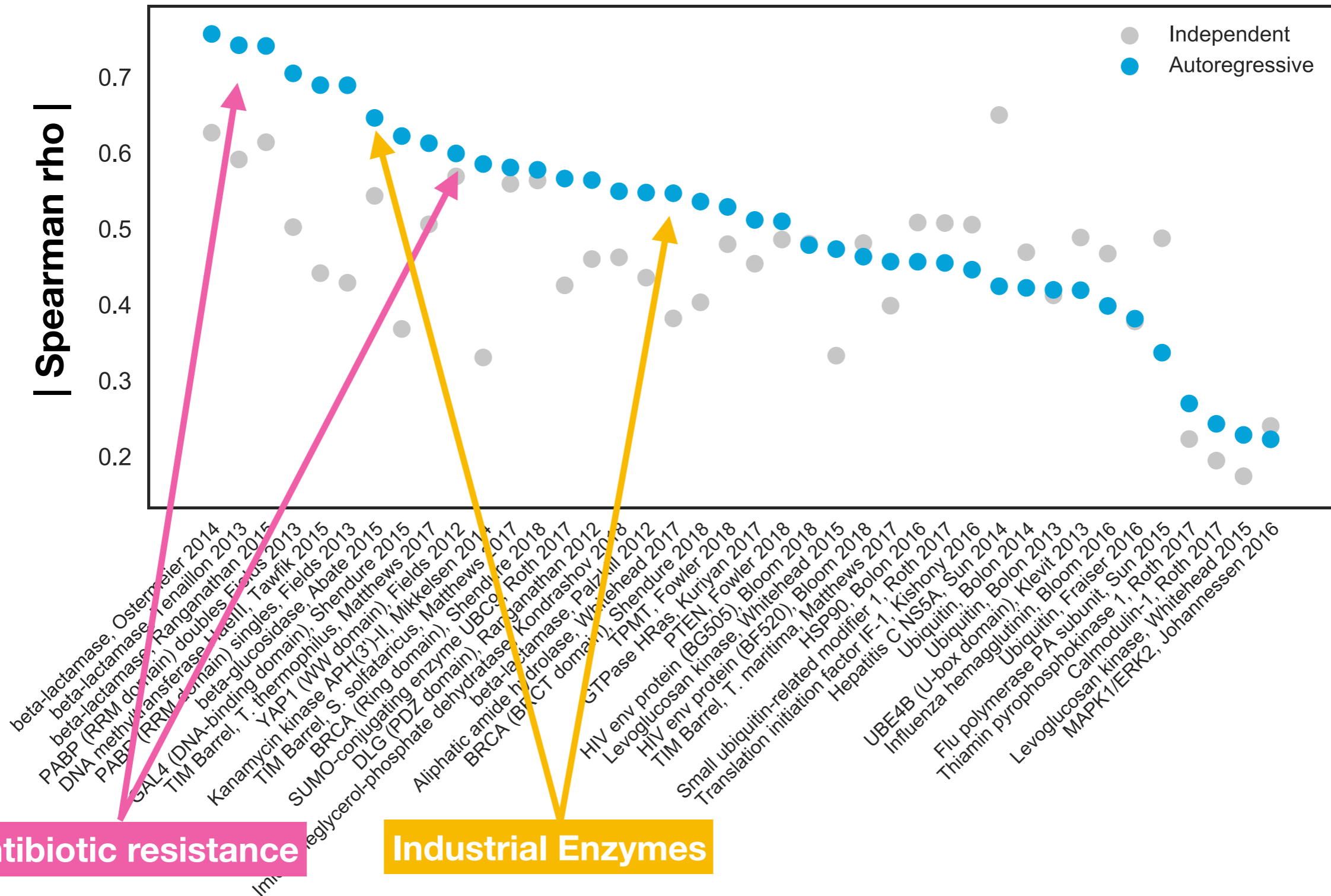
Autoregressive model is predictive of the effect of mutations



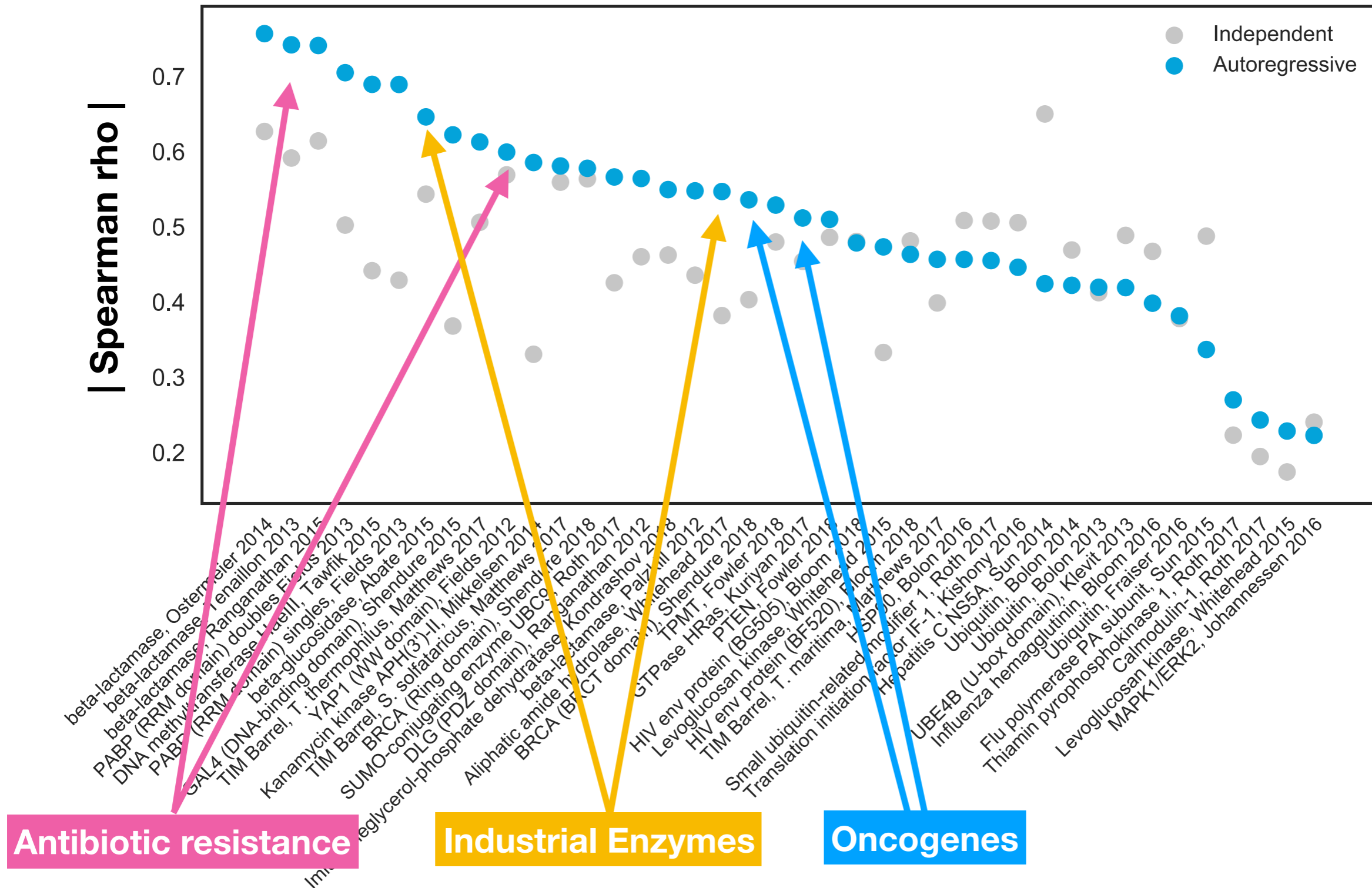
Autoregressive model is predictive of the effect of mutations



Autoregressive model is predictive of the effect of mutations



Autoregressive model is predictive of the effect of mutations



**Autoregressive models are
more flexible**

Autoregressive models are more **flexible**

Missense mutations

...DKSGAG**V**RGSRGIIA...

...DKSGAG**I**RGSRGIIA...

Autoregressive models are more flexible

Missense mutations

...DKSGAGV RGS RGI IA...

...DKSGAG I RGS RGI IA...

Insertions

...DKSGAGEGTRGS RGI IA...

Autoregressive models are more flexible

Missense mutations

...DKSGAGV RGSRGIIA...

...DKSGAGI RGSRGIIA...

Insertions

...DKSGAGEGTRGSRGIIA...

Deletions

...DKSGAGRGSRGIIA...

Models **predict** the **effects** of **insertions** and **deletions**

Insertions

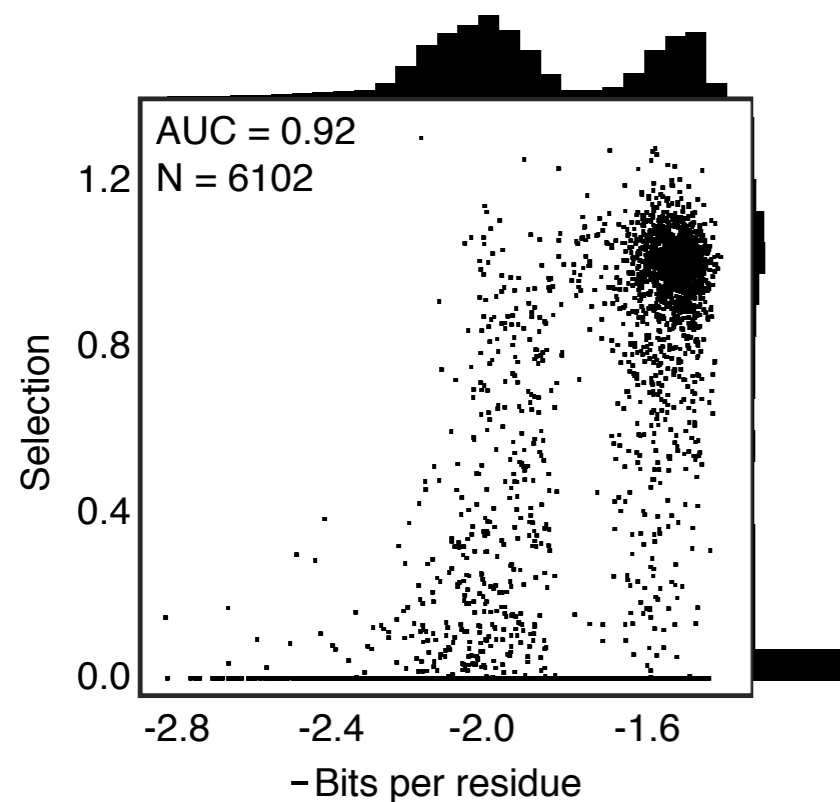
...DKSGAGE**EGT**RGSRGIIA...

Deletions

...DKSGAGRGSRGIIA...

Imidazoleglycerol-phosphate dehydratase

Insertions & deletions



Models predict the effects of insertions and deletions

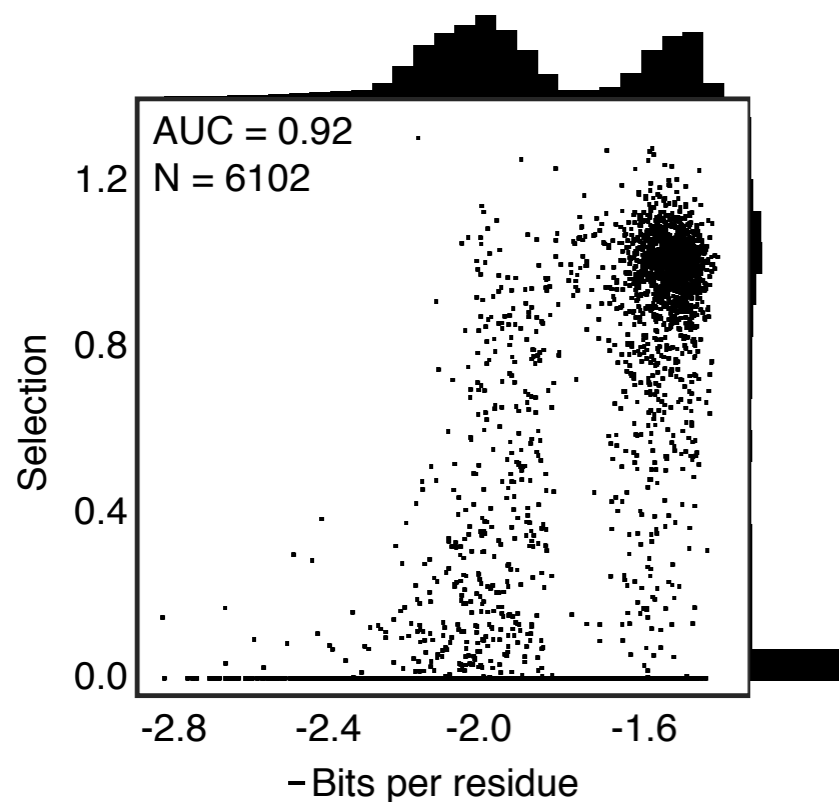
Insertions

...DKSGAGE**EGT**RGSRGIIA...

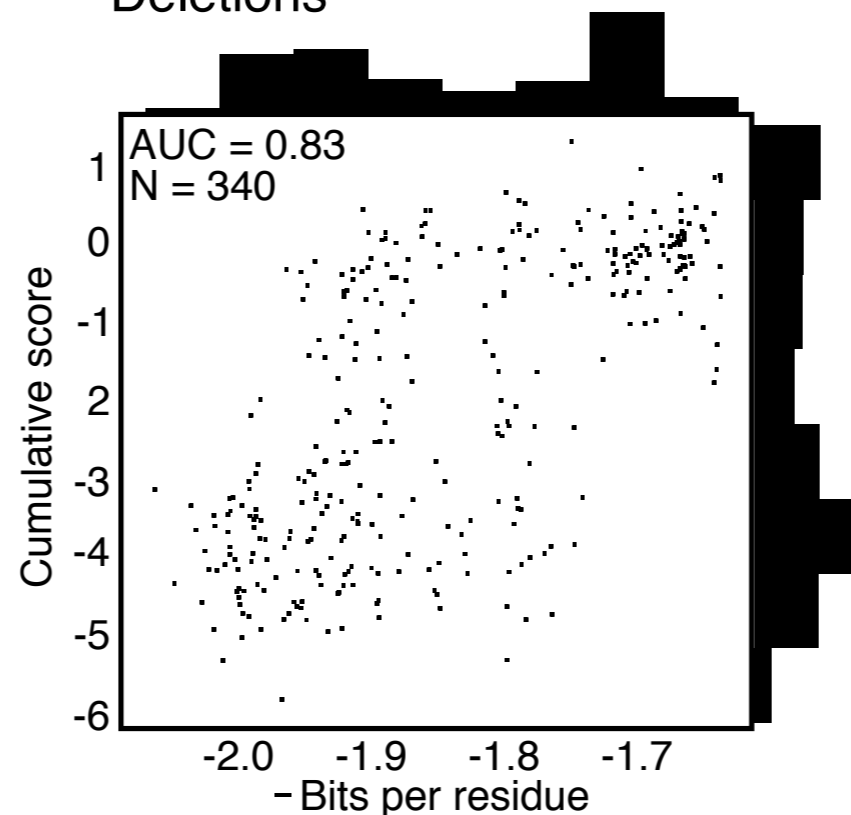
Deletions

...DKSGAGRGSRGIIA...

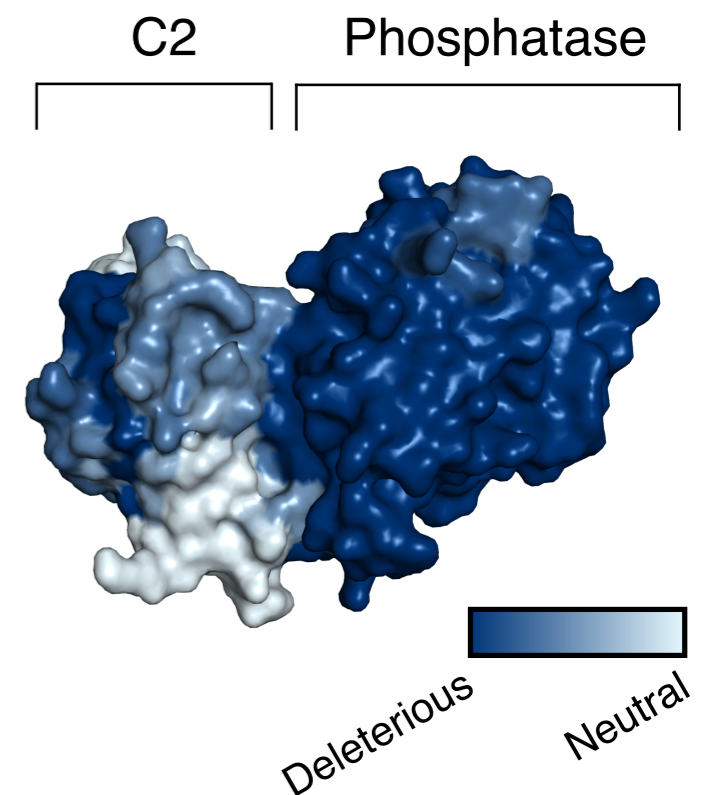
Imidazoglycerol-phosphate dehydratase
Insertions & deletions



PTEN phosphatase
Deletions



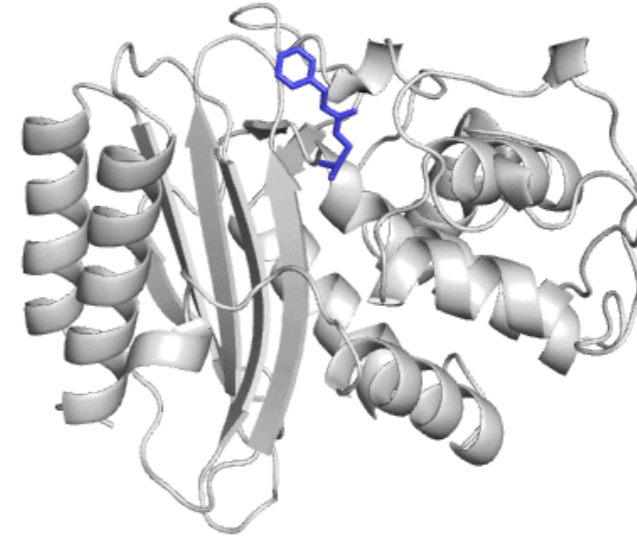
Single amino acid deletions



Recap

Recap

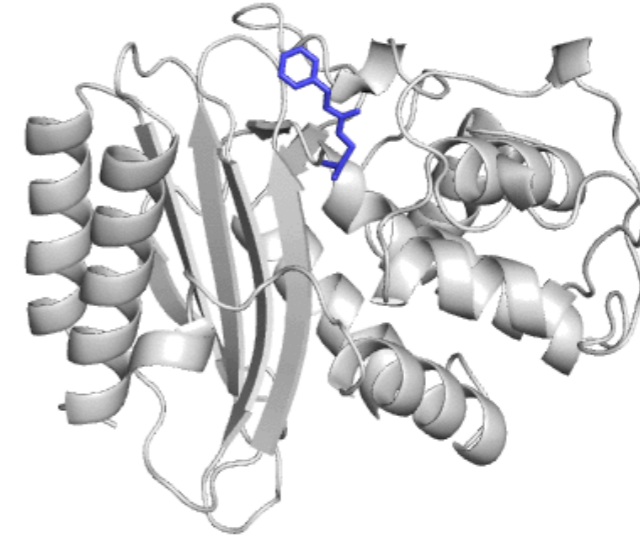
Predicting the effects of mutations is important



Recap

Predicting the effects of mutations is important

Sequencing is becoming cheaper and easier



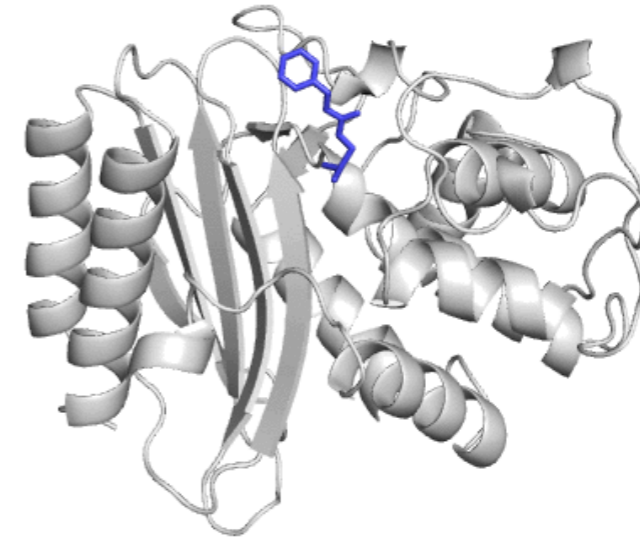
```
GTGTTGCCAACTCGAAGGCTCTGCTCACCAATGTACATGG  
AAAGAGCCGTTTAAATCTCGGGCGCTTAATTCGCCTCGTGA  
TTCCAGCCAGGCAGTGGCGGATCAATATGCCGACTTCCTG  
ACACCTCGTCGATGGATTACTACCATCAGTTGCGTTATGC  
GTTGGGGCTGGATTACCGGTTATTGAGAACCTGCAAAATC  
TCTTTCTGGTTCGCTTTCTTATATCTTCGGCAAGTTAGAC  
AAATGGGTTATACCGAACCGGACCCGCGAGATGATCTTTC  
GAAACGGGACGTGAACTGGAGCTGGCGGATATTGAAATTG  
TGCCGCTTTTATGGCGAATCTGTCACAACCTCGACGATCTC
```


Recap

Predicting the effects of mutations is important

Sequencing is becoming cheaper and easier

Generative models fit to biological sequence data can predict the effect of mutations



```
GTGTTGCCAACTCGAAGGCTCTGCTCACCAATGTACATGG  
AAAGAGCCGTTTAAATCTCGGGCGCTTAATTCGCCTCGTGA  
TTCCAGCCAGGCAGTGGCGGATCAATATGCCGACTTCCTG  
ACACCTCGTCGATGGATTACTACCATCAGTTGCGTTATGC  
GTTGGGGCTGGATTACCGGTTATTGAGAACCTGCAAAATC  
TCTTTCTGGTTCGCTTTCTTATATCTTCGGCAAGTTAGAC  
AAATGGGTTATACCGAACCGGACCCGCGAGATGATCTTTC  
GAAACGGGACGTGAACTGGAGCTGGCGGATATTGAAATTG  
TGCCGCTTTTATGGCGAATCTGTCACAACCTCGACGATCTC
```

$$\log \frac{p(\mathbf{x}_{\text{mut}} | \boldsymbol{\theta})}{p(\mathbf{x}_{\text{wild}} | \boldsymbol{\theta})}$$

Part II: **Genotype -> phenotype** for **complex diseases**



Part II: **Genotype -> phenotype** for **complex diseases**



Part II: **Genotype -> phenotype** for **complex diseases**



Founded by Daphne Koller in 2018

Part II: **Genotype -> phenotype** for **complex diseases**



Founded by Daphne Koller in 2018

Create a **new paradigm** for **drug development**

Part II: **Genotype -> phenotype** for **complex diseases**



Founded by Daphne Koller in 2018

Create a **new paradigm** for **drug development**
that uses **high-quality data** and **data-driven models**

Part II: **Genotype -> phenotype** for **complex diseases**



Founded by Daphne Koller in 2018

Create a **new paradigm** for **drug development**
that uses **high-quality data** and **data-driven models**
to design **novel, safe, and effective therapies**

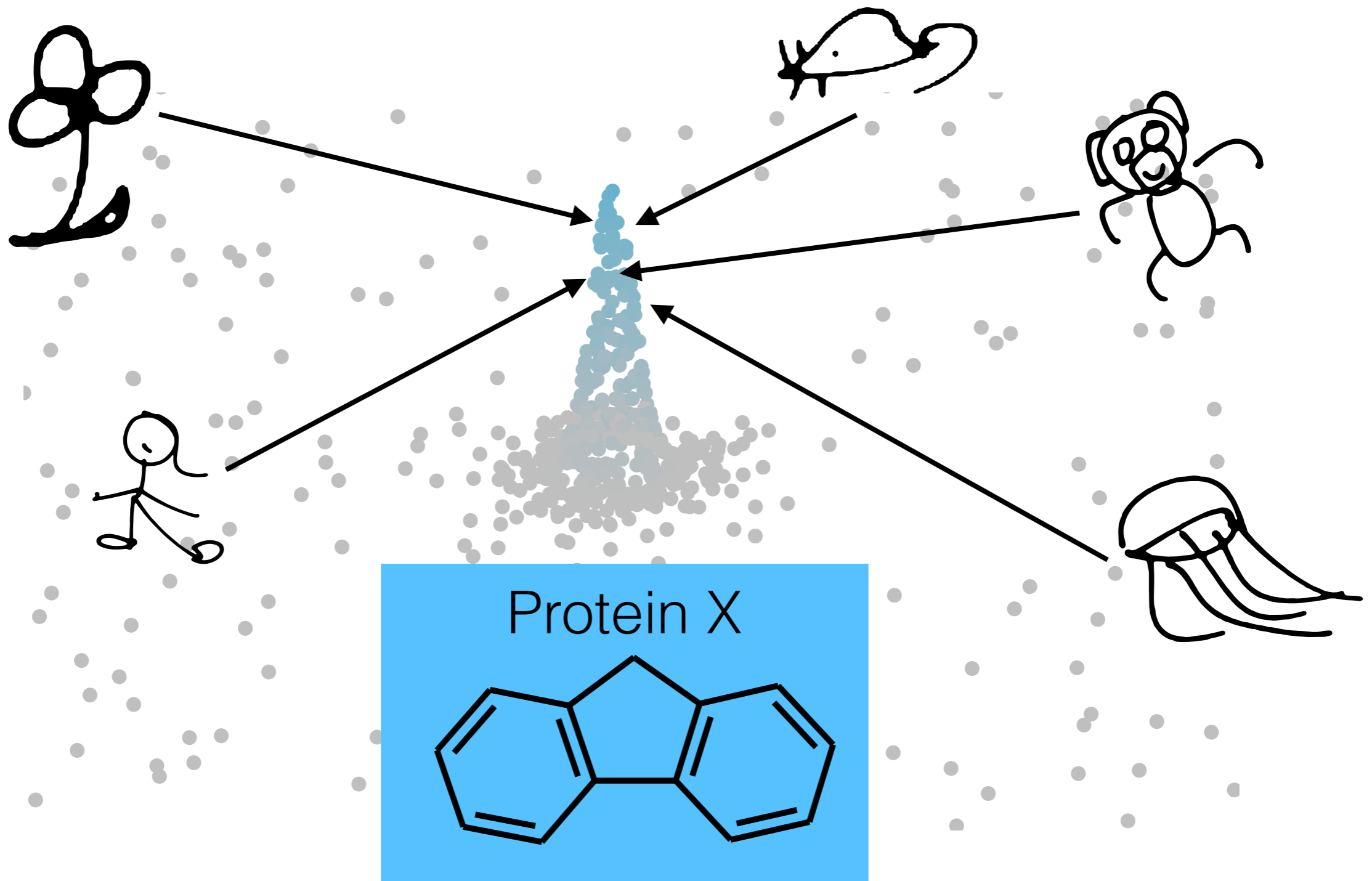
Part II: **Genotype -> phenotype** for **complex diseases**



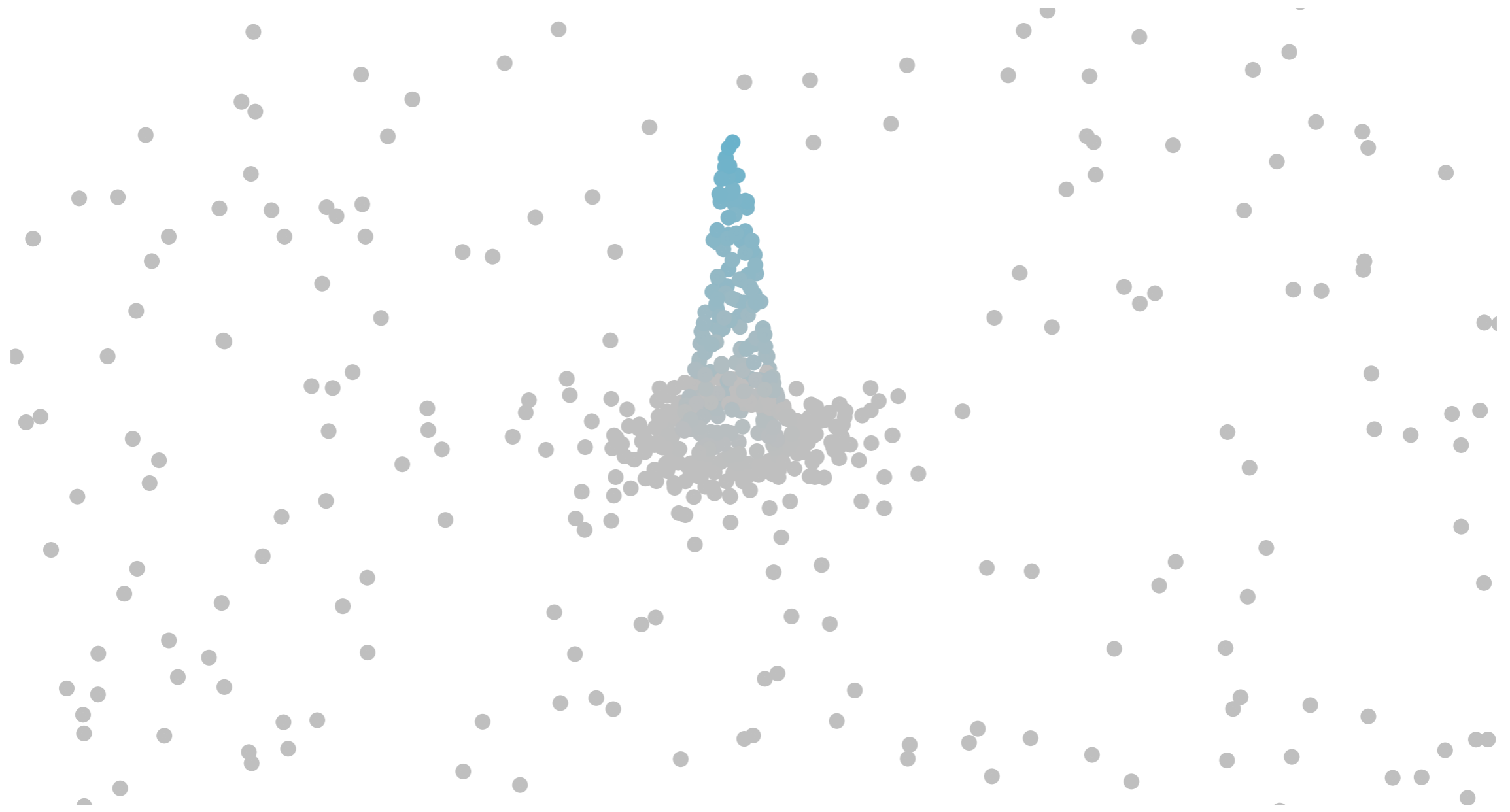
Founded by Daphne Koller in 2018

Create a **new paradigm** for **drug development**
that uses **high-quality data** and **data-driven models**
to design **novel, safe, and effective therapies**
that help **more people, faster, and at a lower cost.**

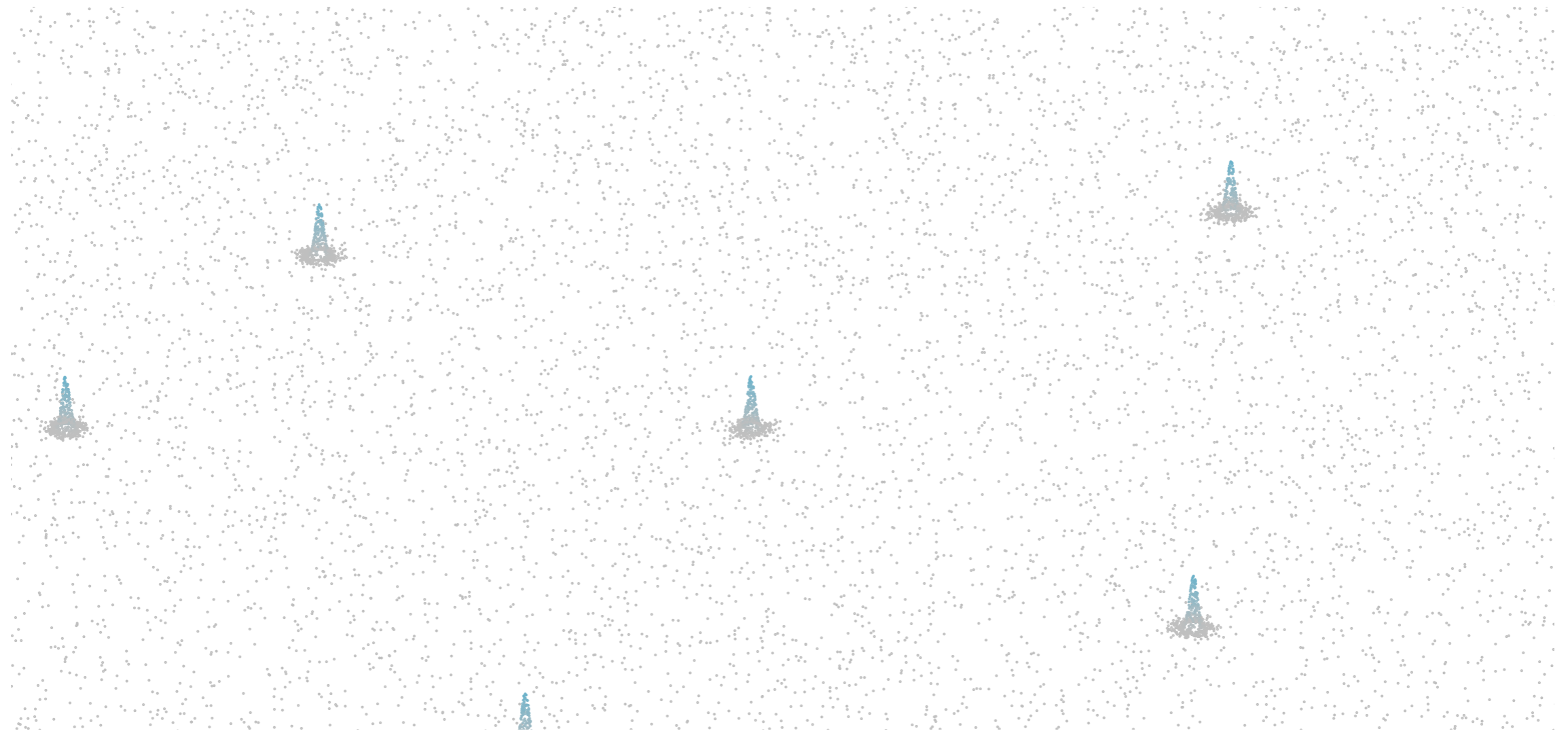
Part II: Genotype -> phenotype for different **tissues** are **sparse**



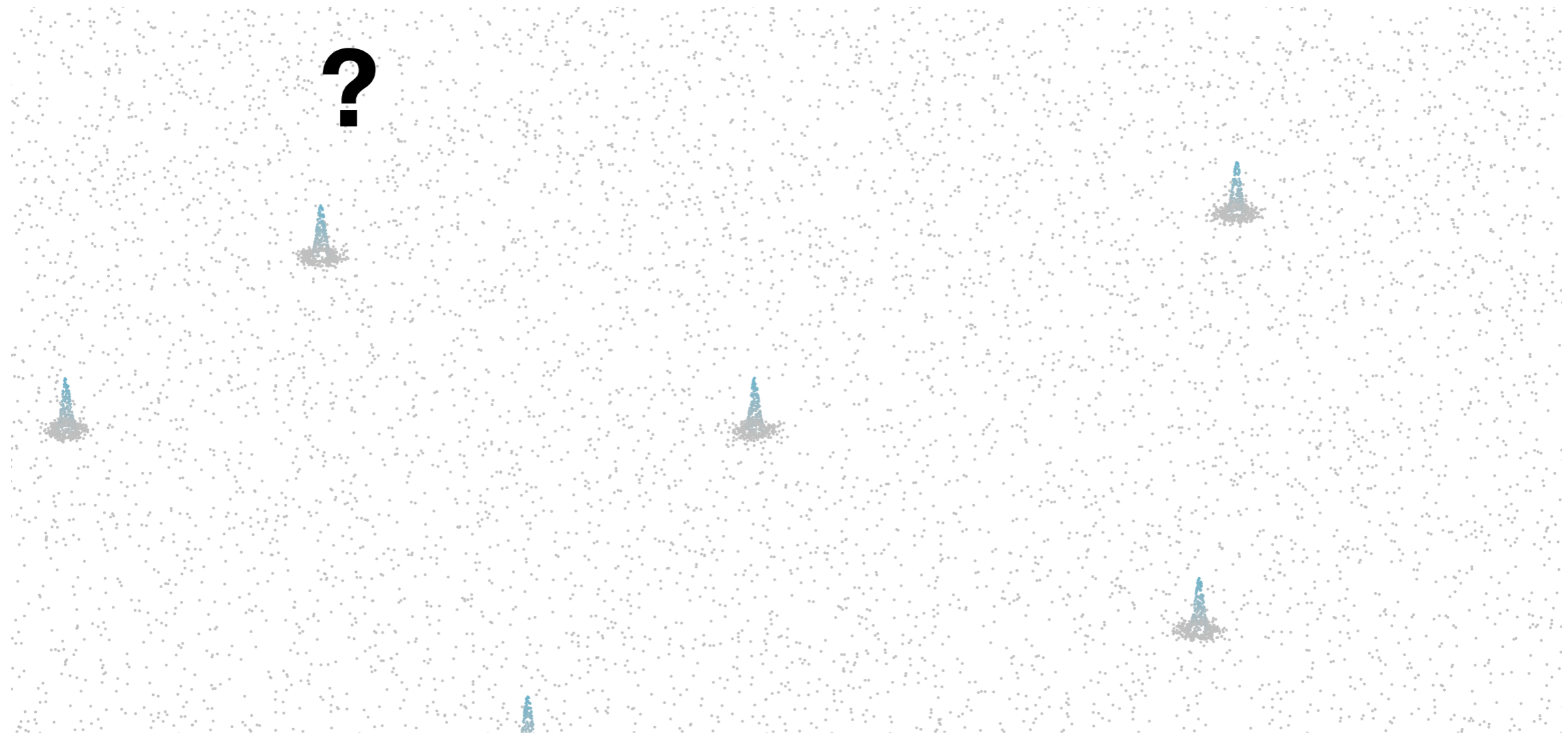
Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**

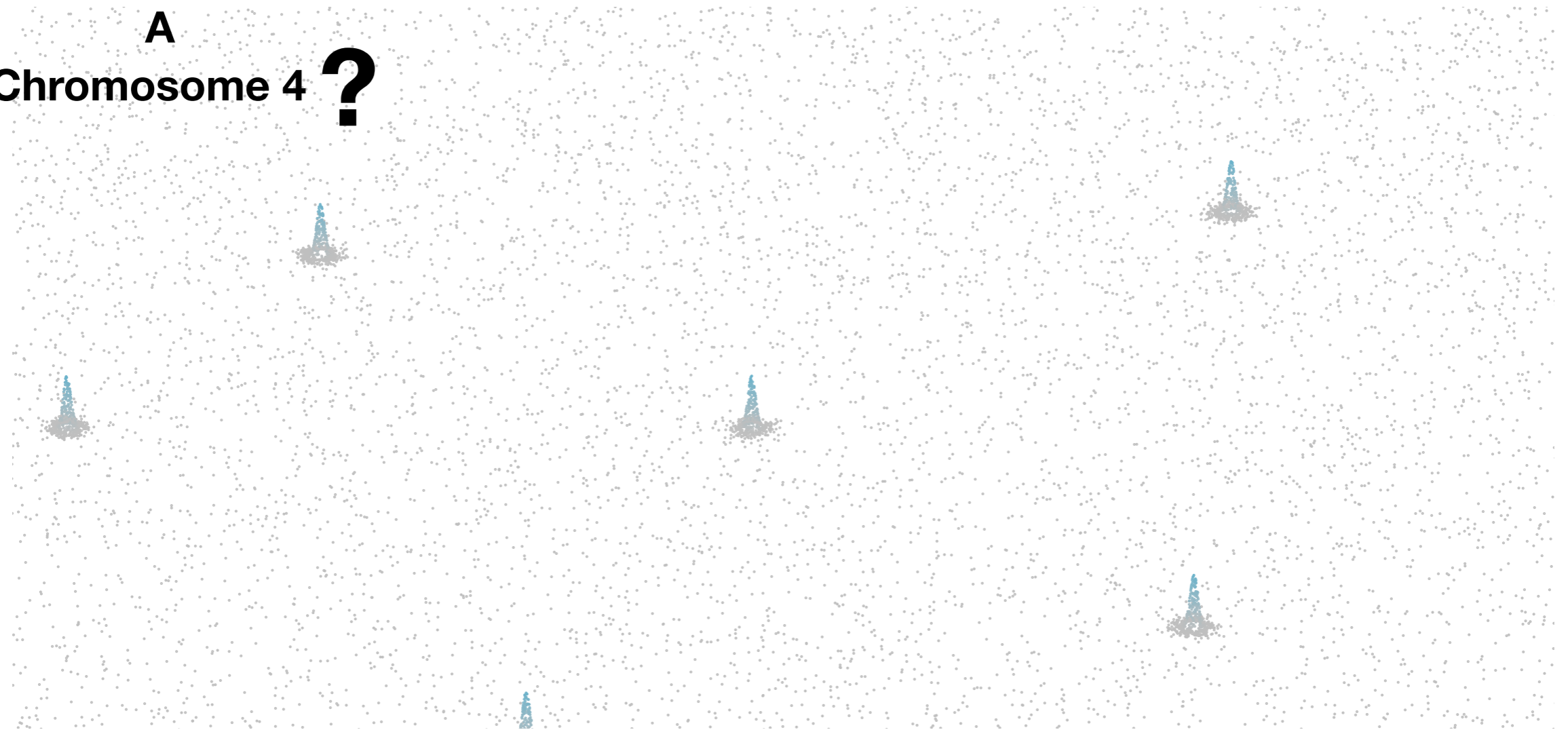


Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**

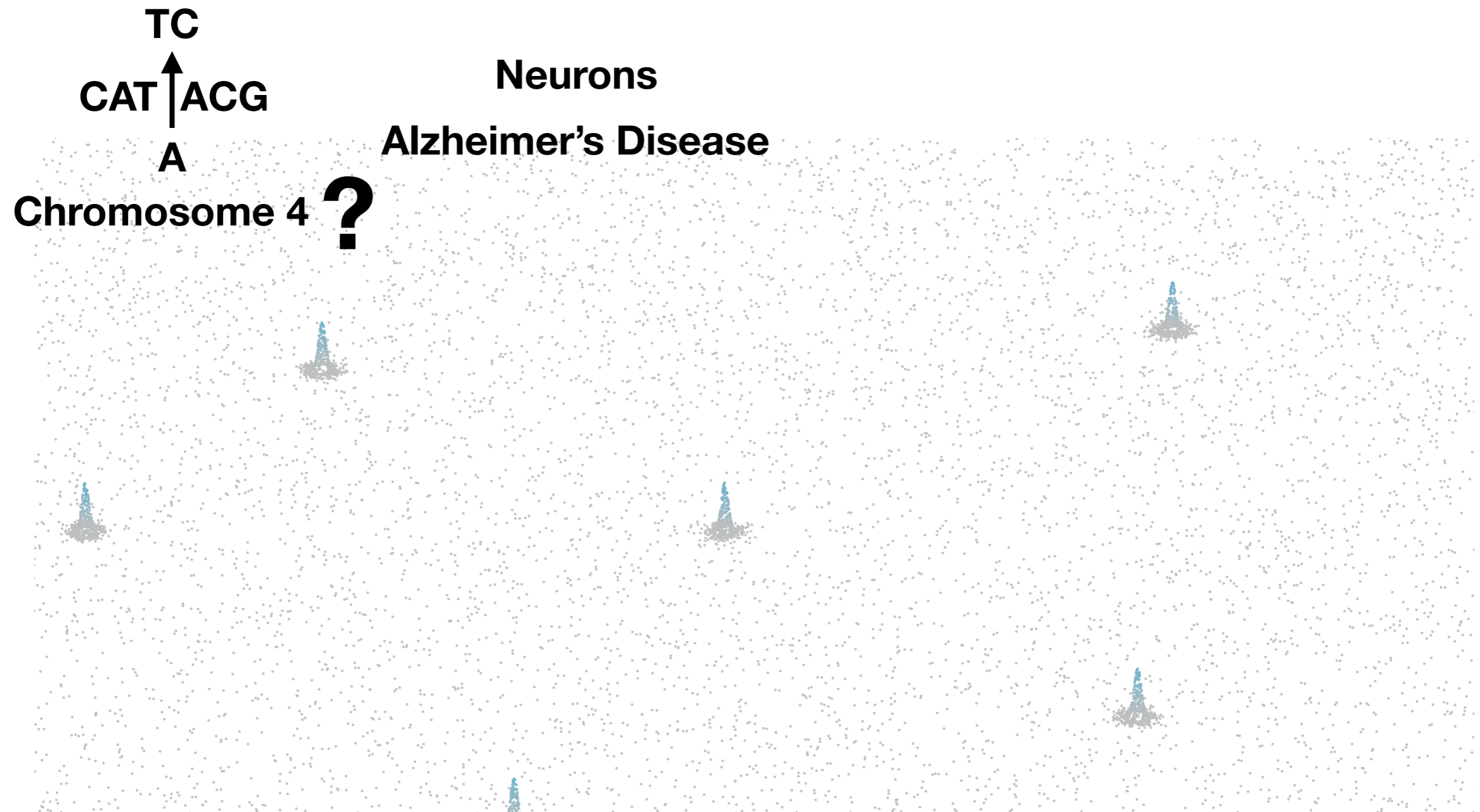


Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**

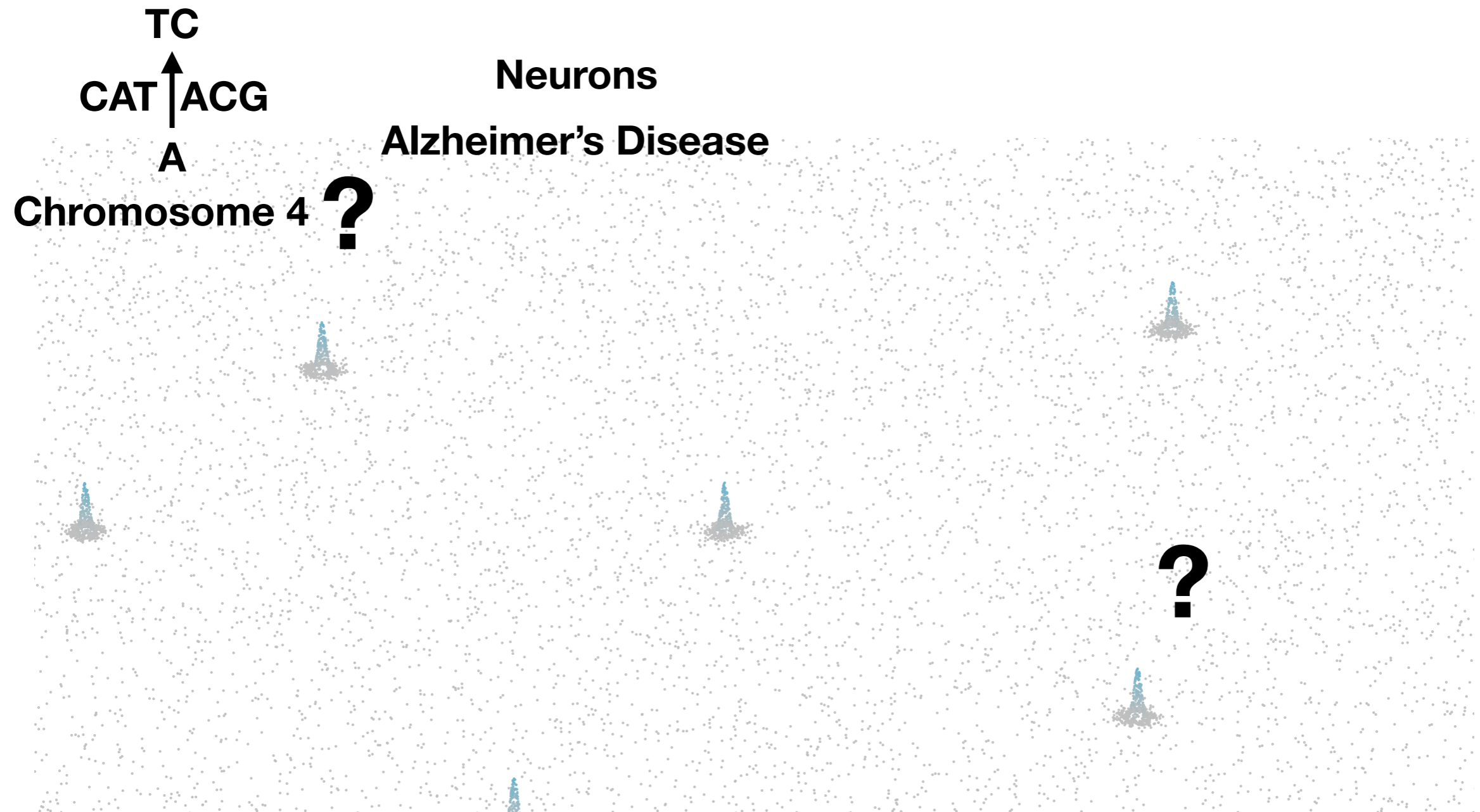
TC
CAT | ACG
A
Chromosome 4 ?



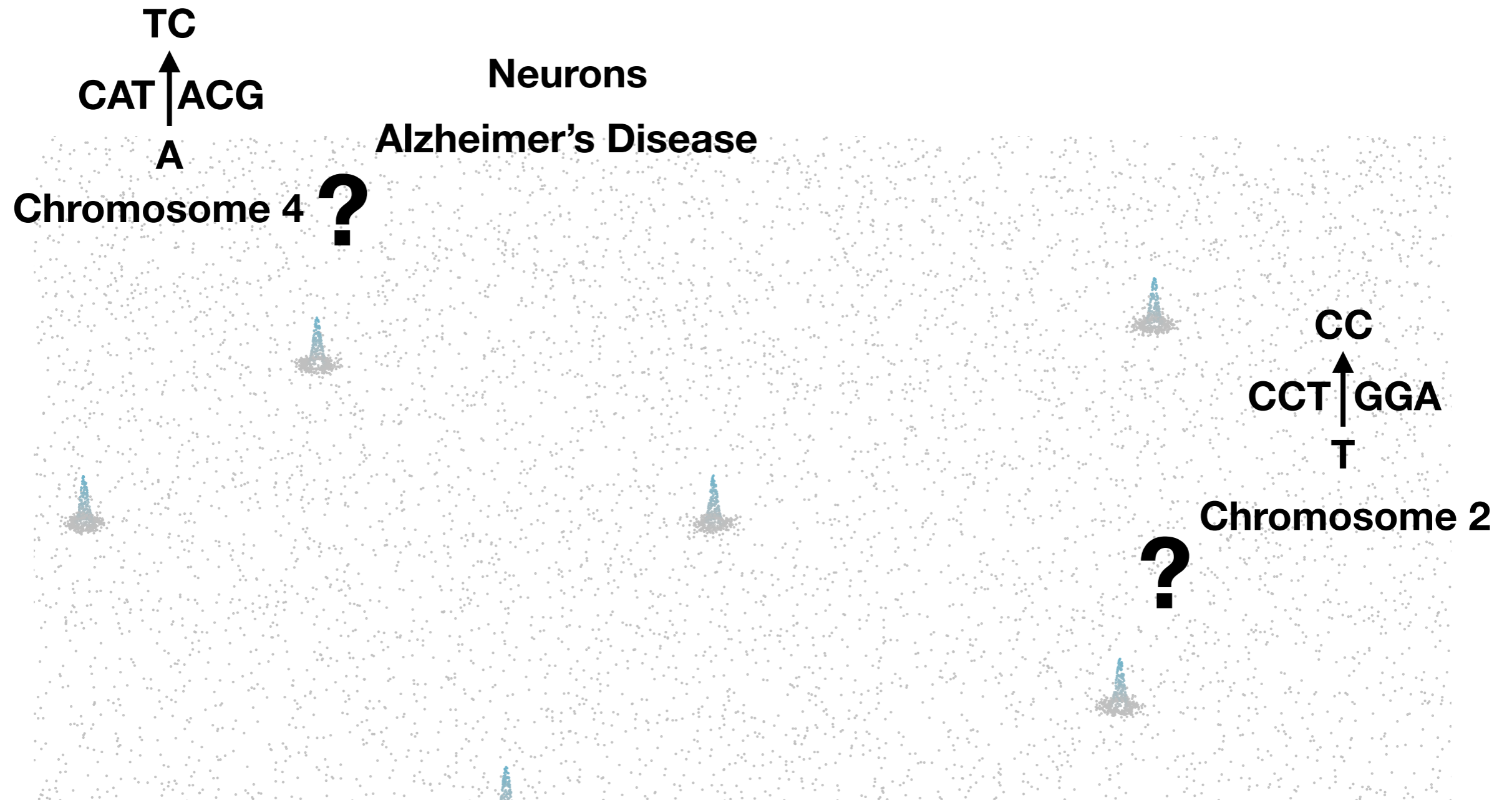
Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



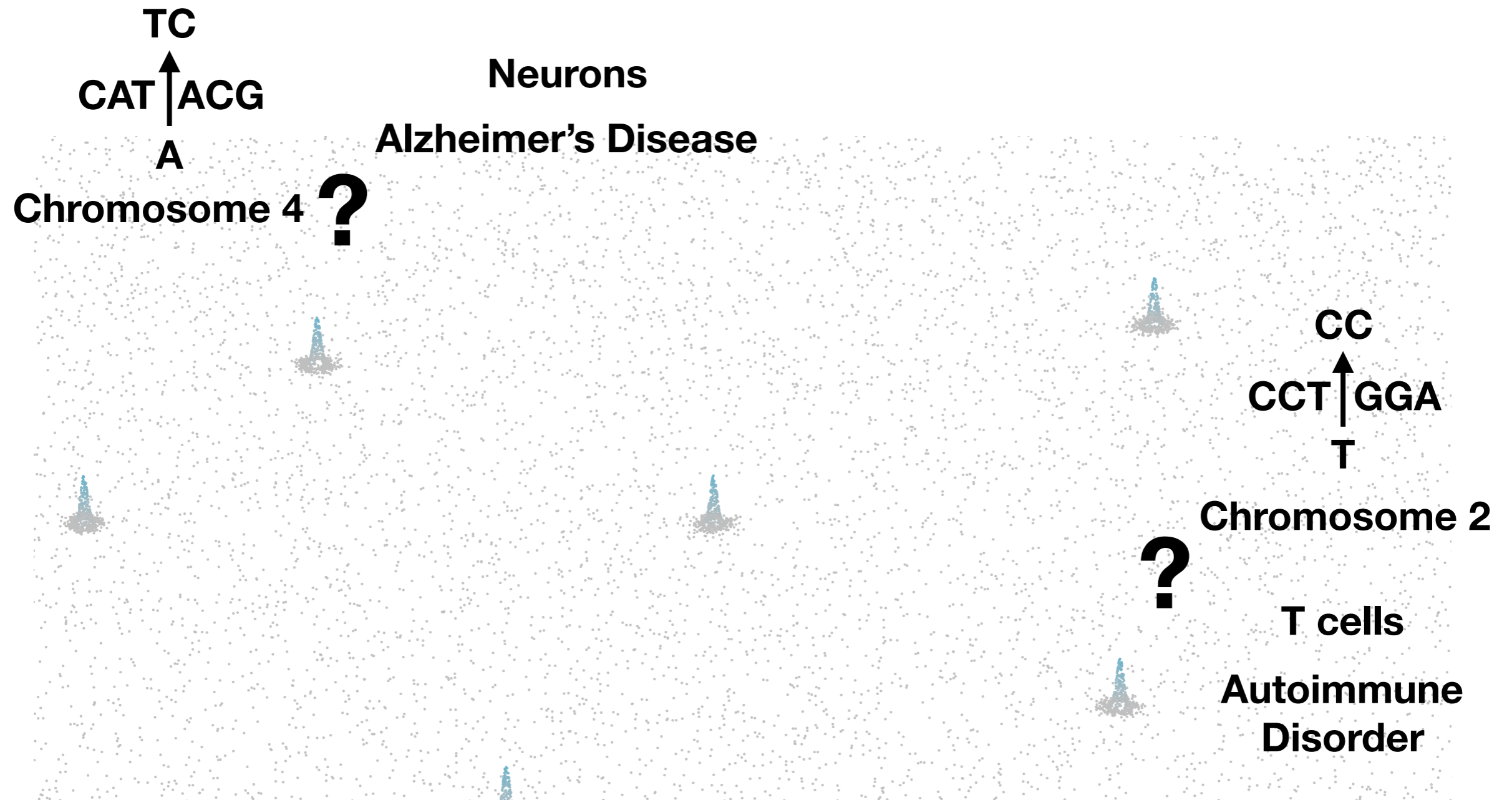
Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



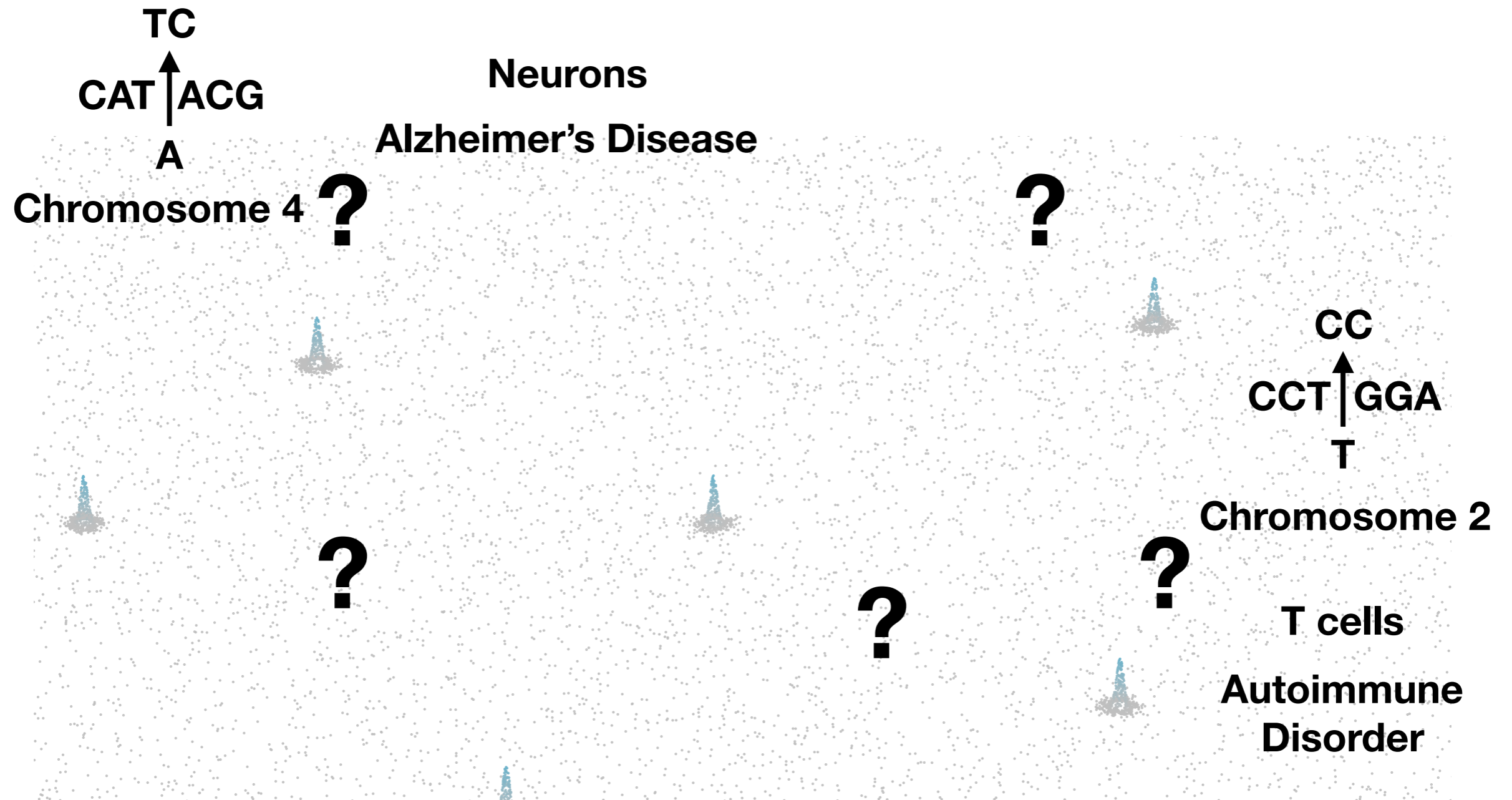
Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



Part II: **Genotype** -> **phenotype** for different **tissues** are **sparse**



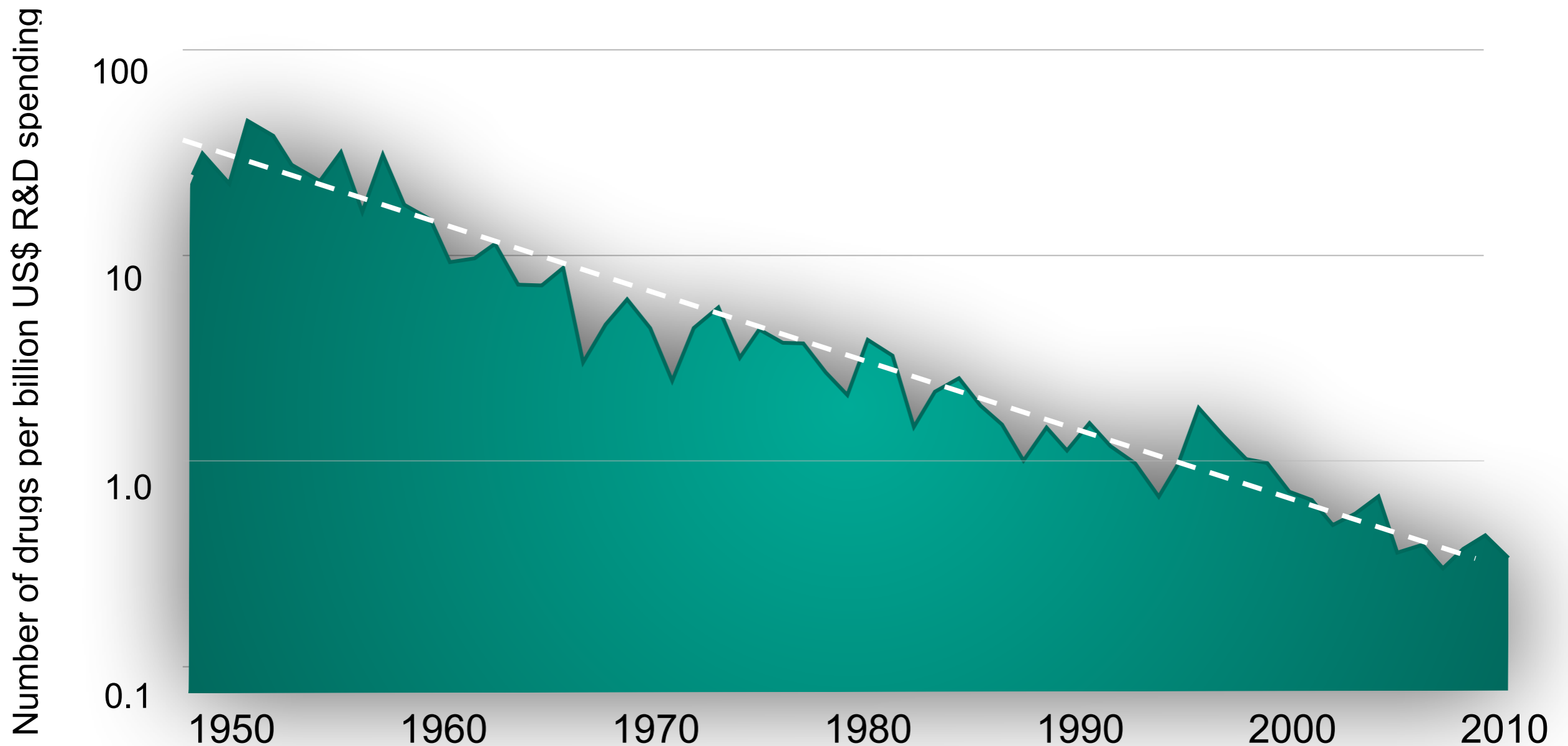
Drug discovery is becoming
exponentially **less efficient**

Drug discovery is becoming
exponentially **less efficient**

Moore's Law  **Eroom's Law**

Drug discovery is becoming exponentially **less efficient**

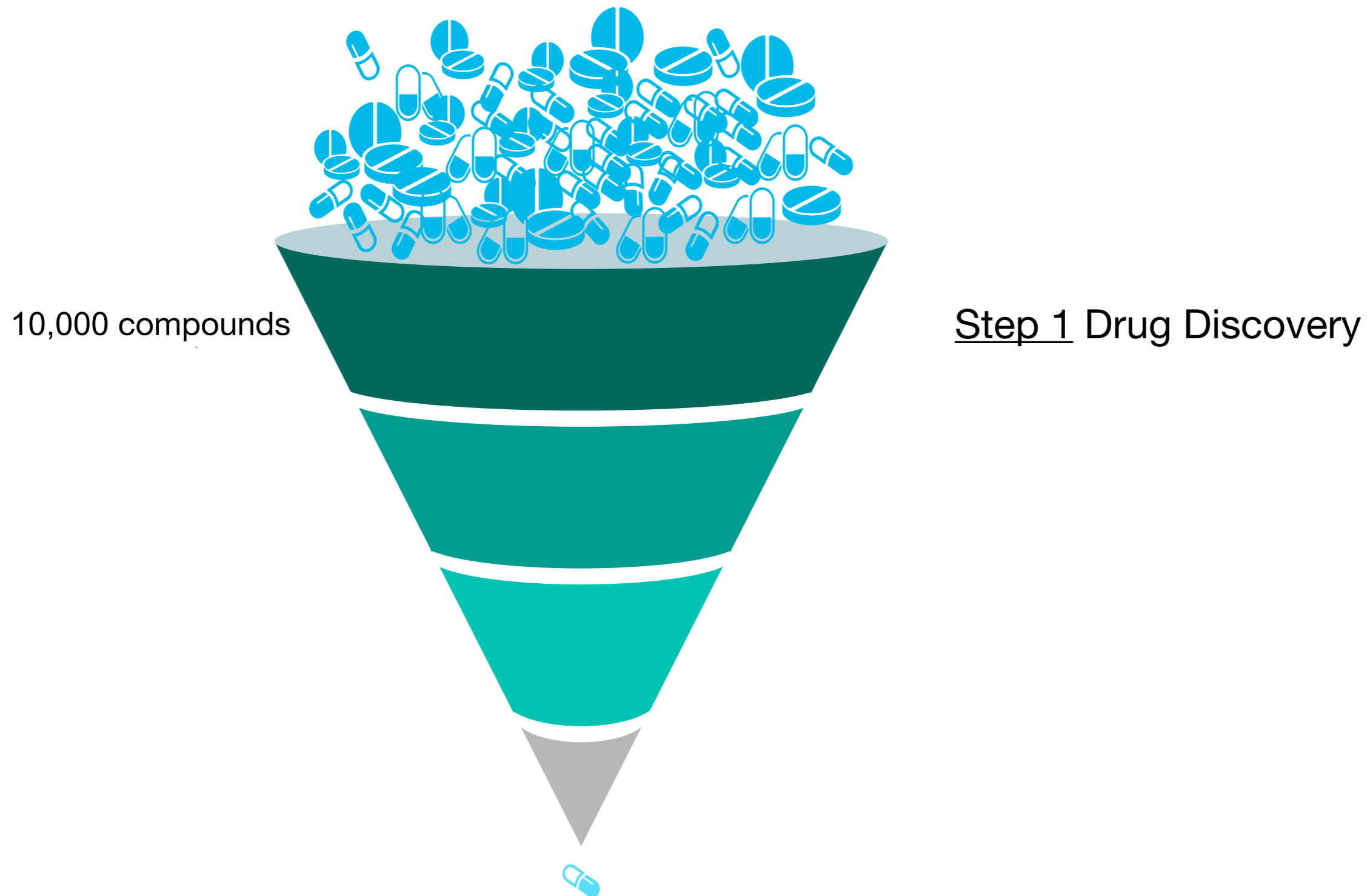
Moore's Law \longrightarrow **Eroom's Law**



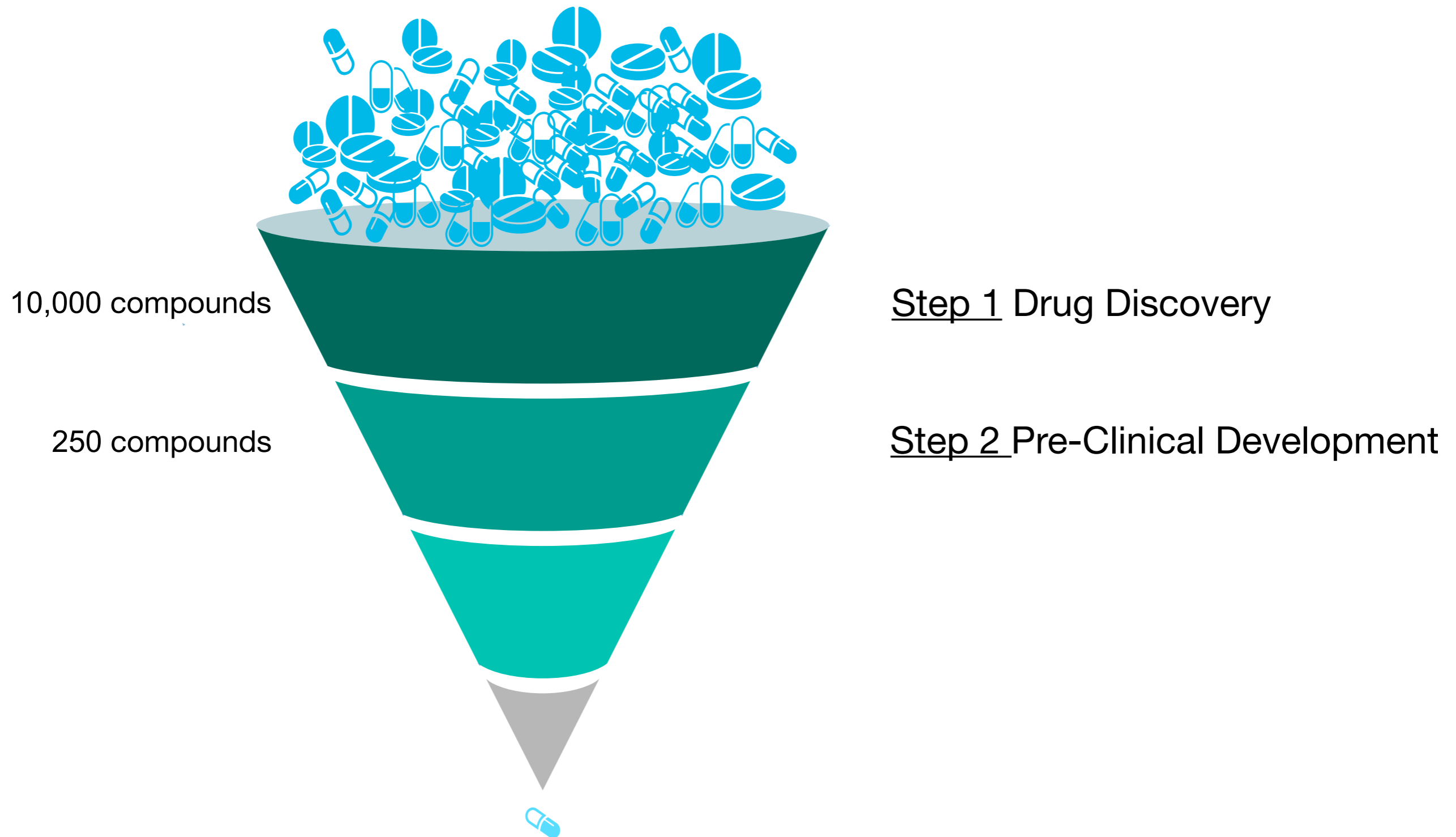
Drug discovery is ultimately
a **problem of prediction**



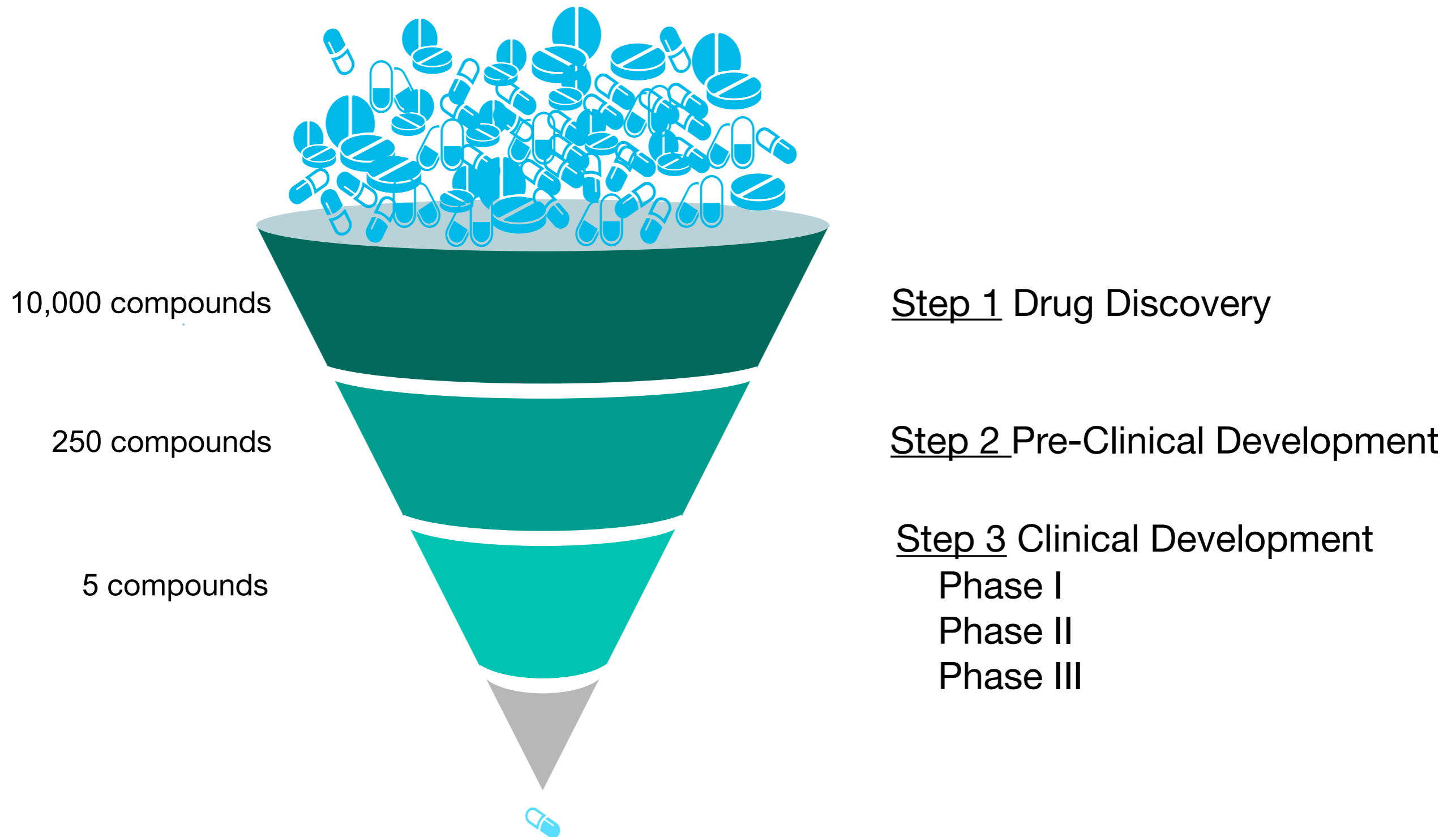
Drug discovery is ultimately a problem of prediction



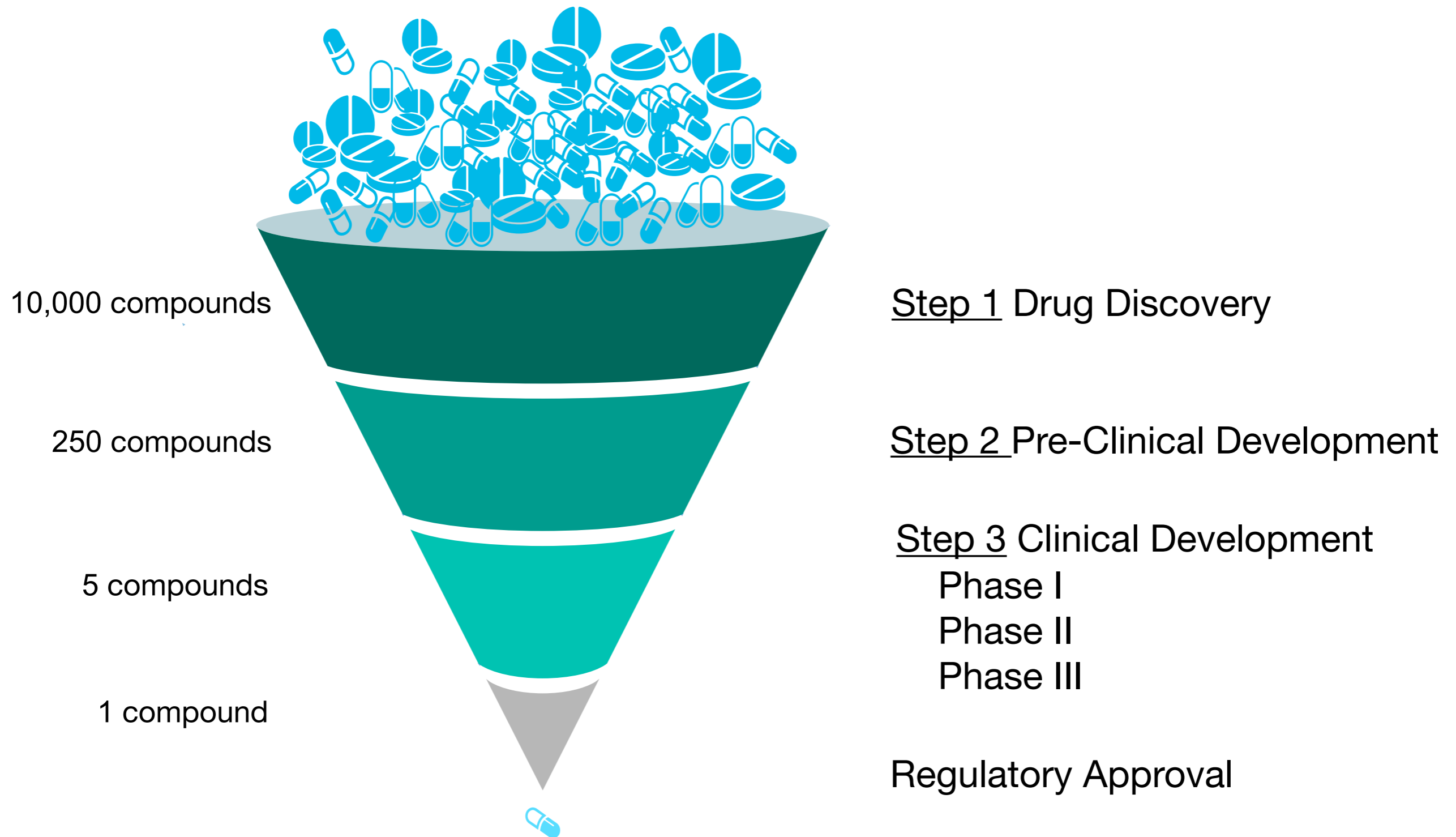
Drug discovery is ultimately a problem of prediction



Drug discovery is ultimately a problem of prediction



Drug discovery is ultimately a problem of prediction



Drug discovery is ultimately a problem of prediction

What can we do to make this pipeline better?



Drug discovery is ultimately a problem of prediction

What can we do to make this pipeline better?

10,000 compounds

Step 1 Drug Discovery

250 compounds

Step 2 Pre-Clinical Development

High-throughput biology

Step 3 Clinical Development

5 compounds

Phase I

+

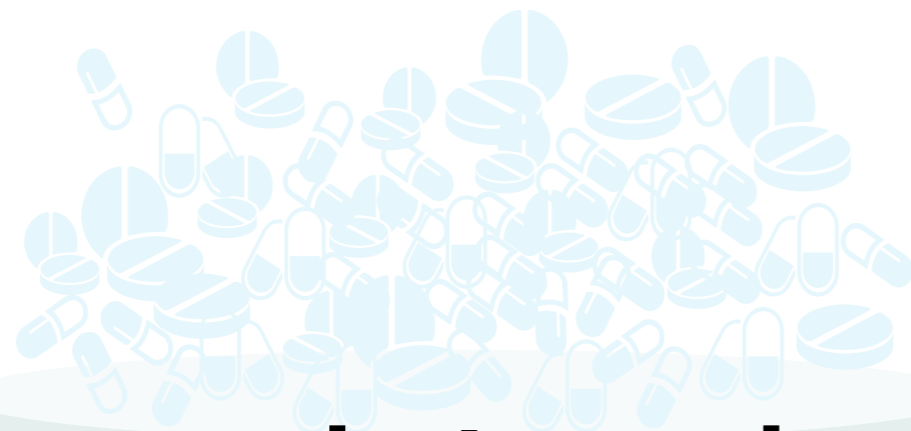
Phase II

Computation

Phase III

1 compound

Regulatory Approval



insitro is founded on **data-driven** technologies

**Human
genetics**



insitro is founded on **data-driven** technologies

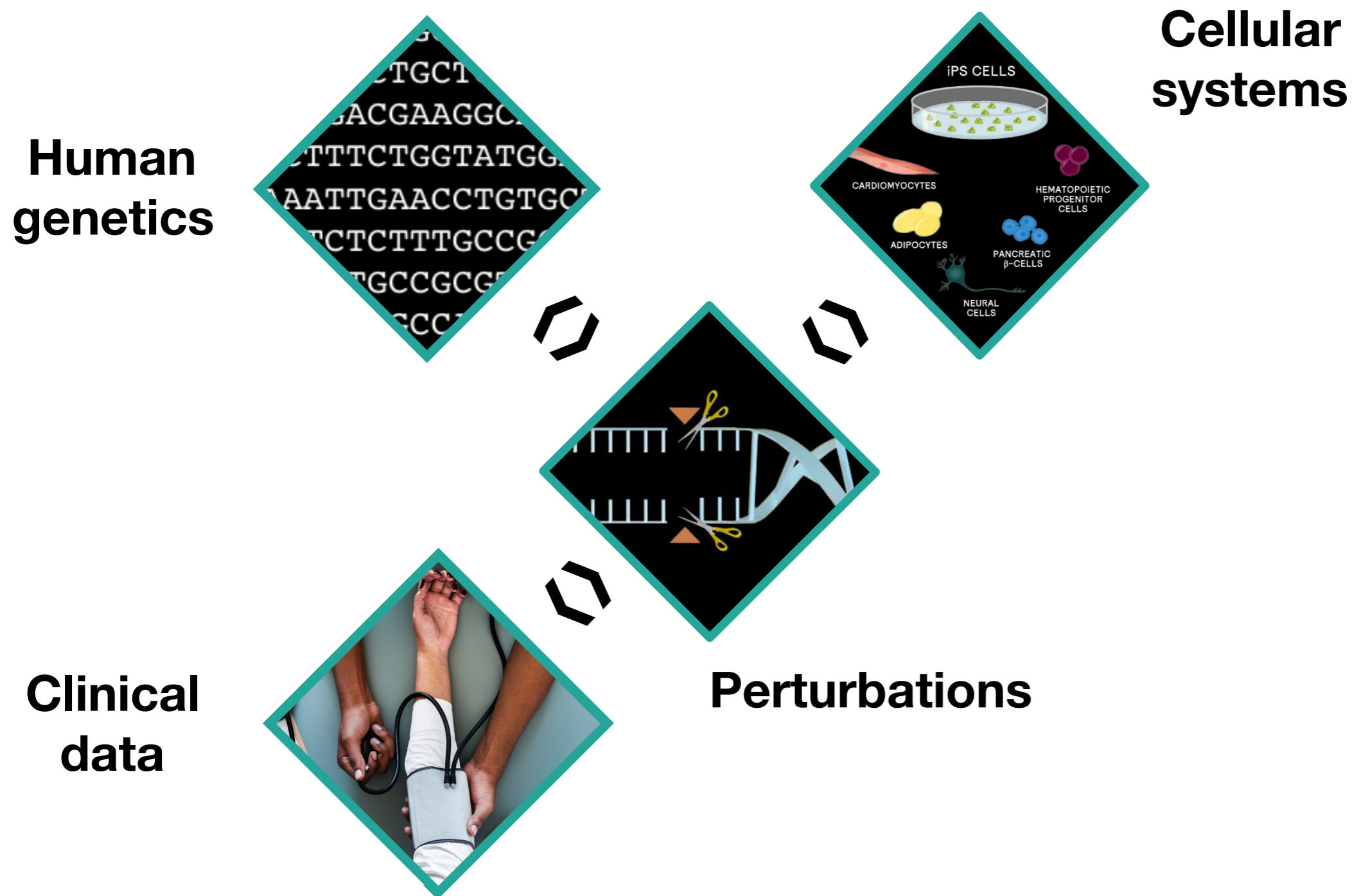
**Human
genetics**



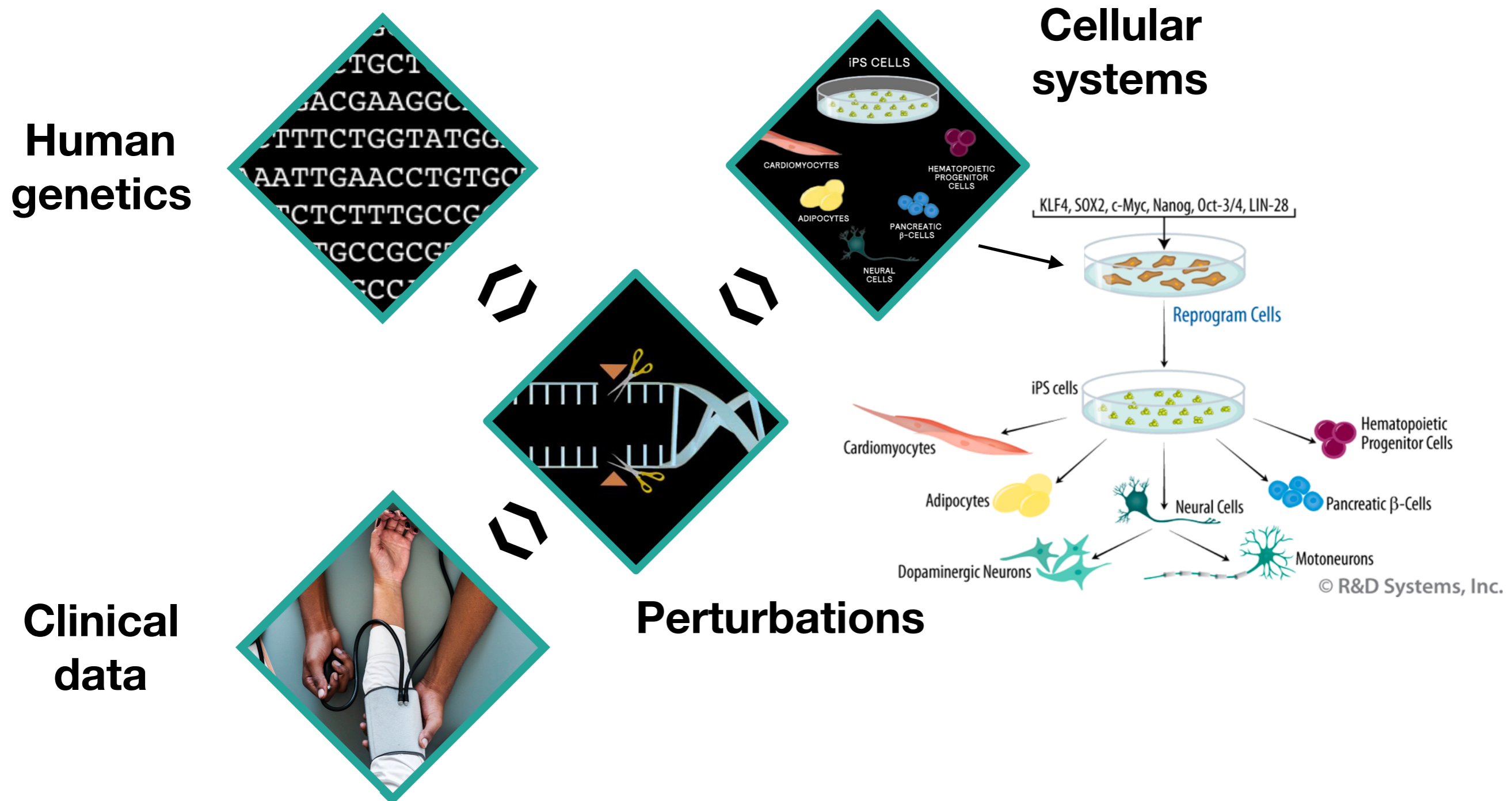
**Clinical
data**



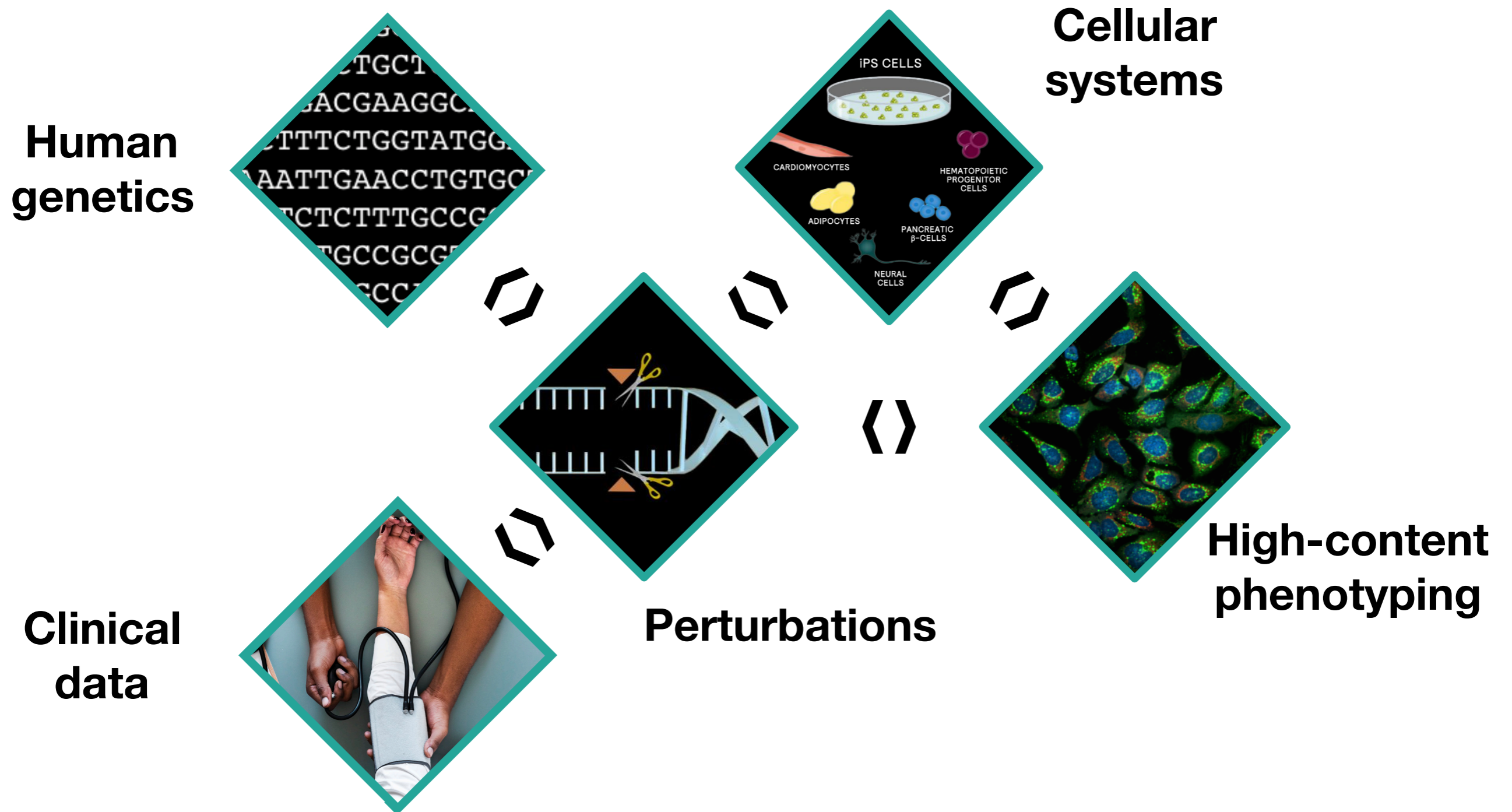
insitro is founded on data-driven technologies



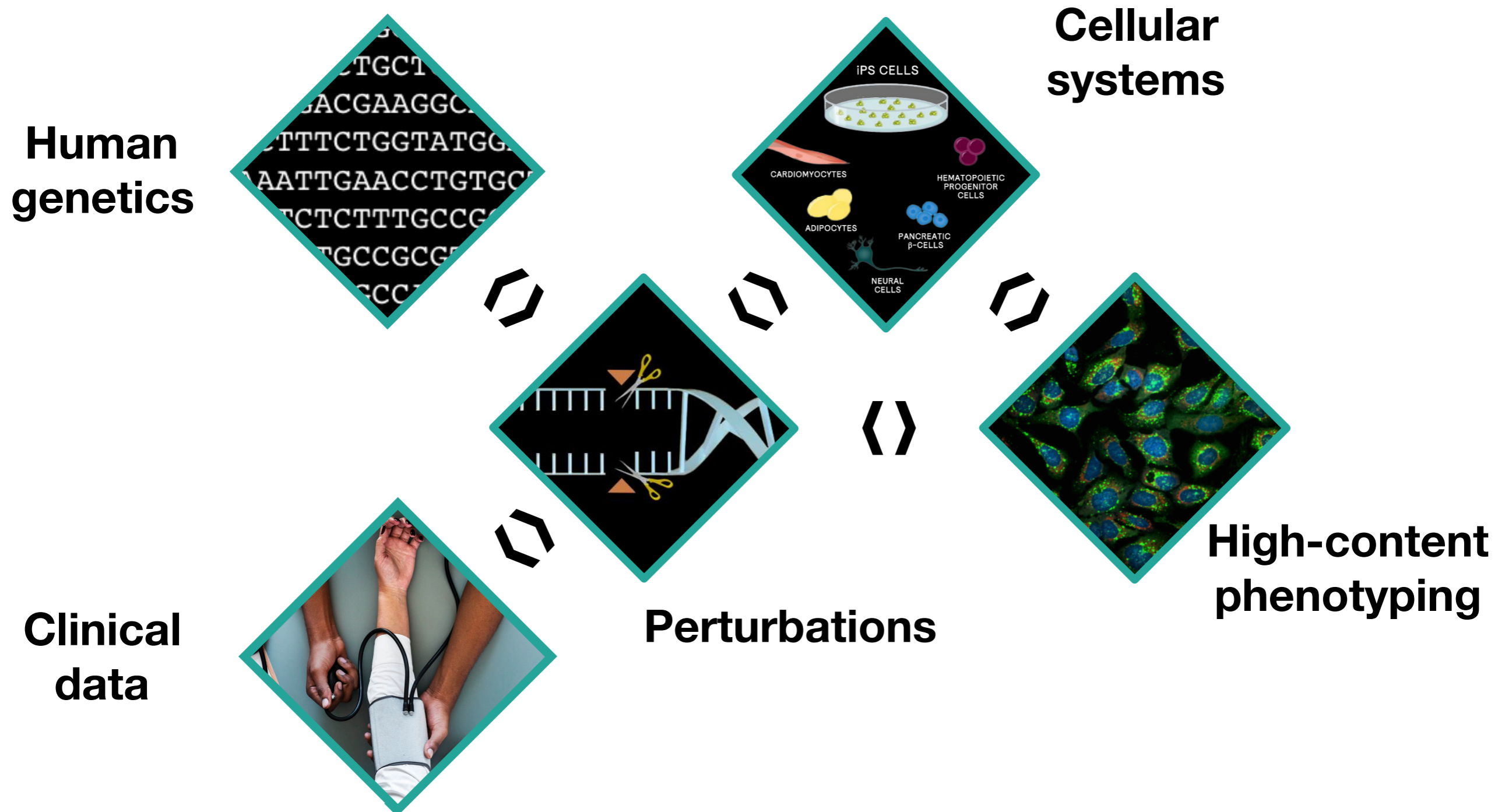
insitro is founded on data-driven technologies



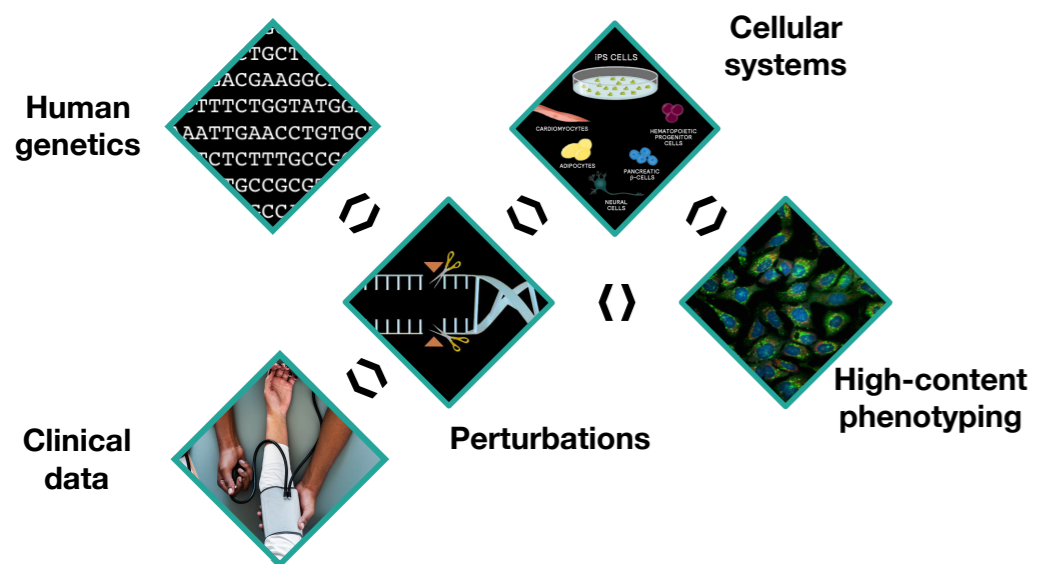
insitro is founded on data-driven technologies



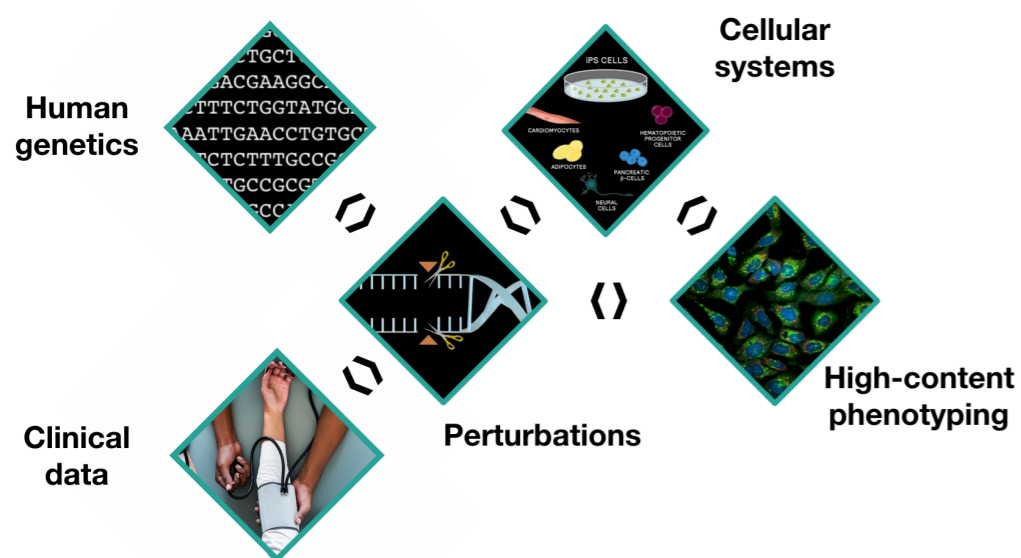
insitro is founded on data-driven technologies



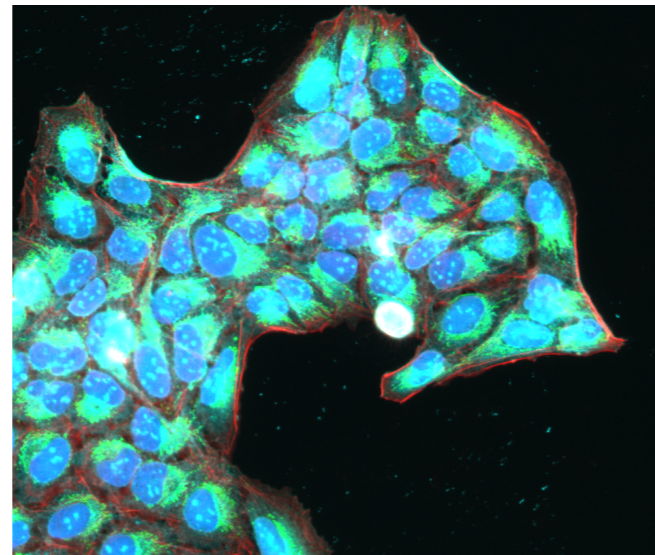
insitro is founded on data-driven technologies



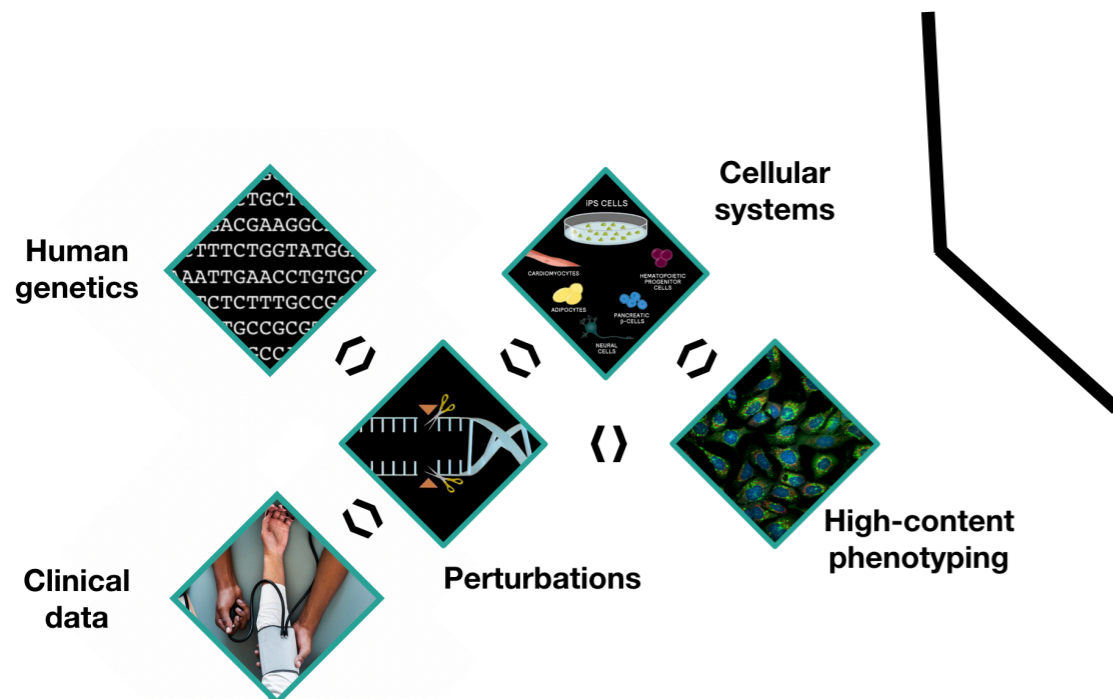
insitro is founded on data-driven technologies



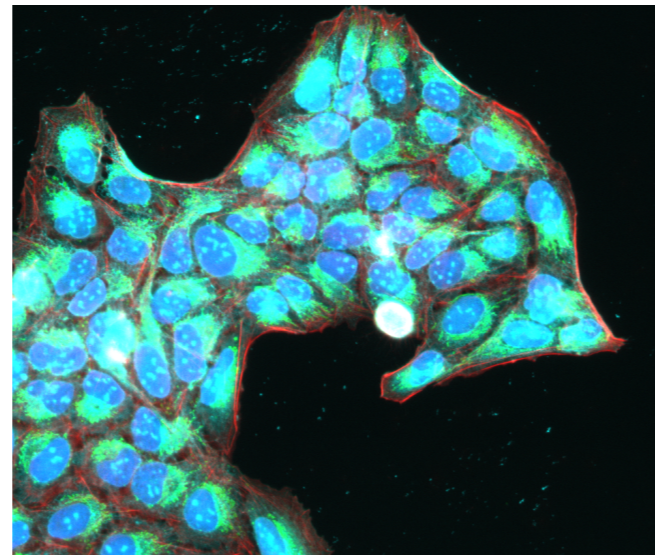
insitro is founded on data-driven technologies



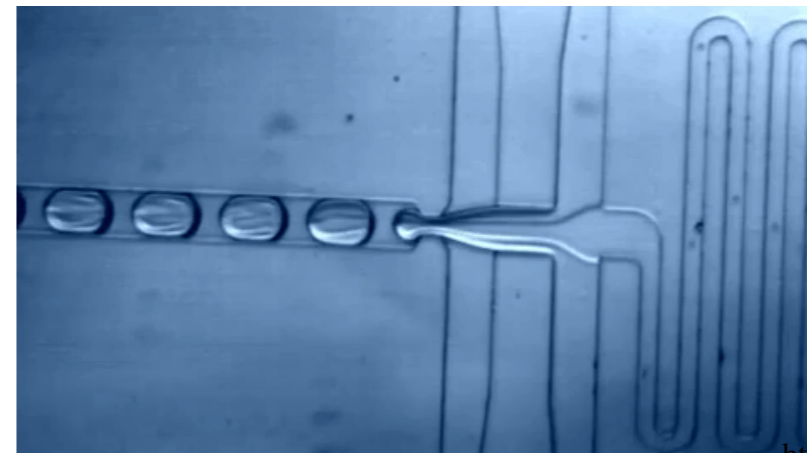
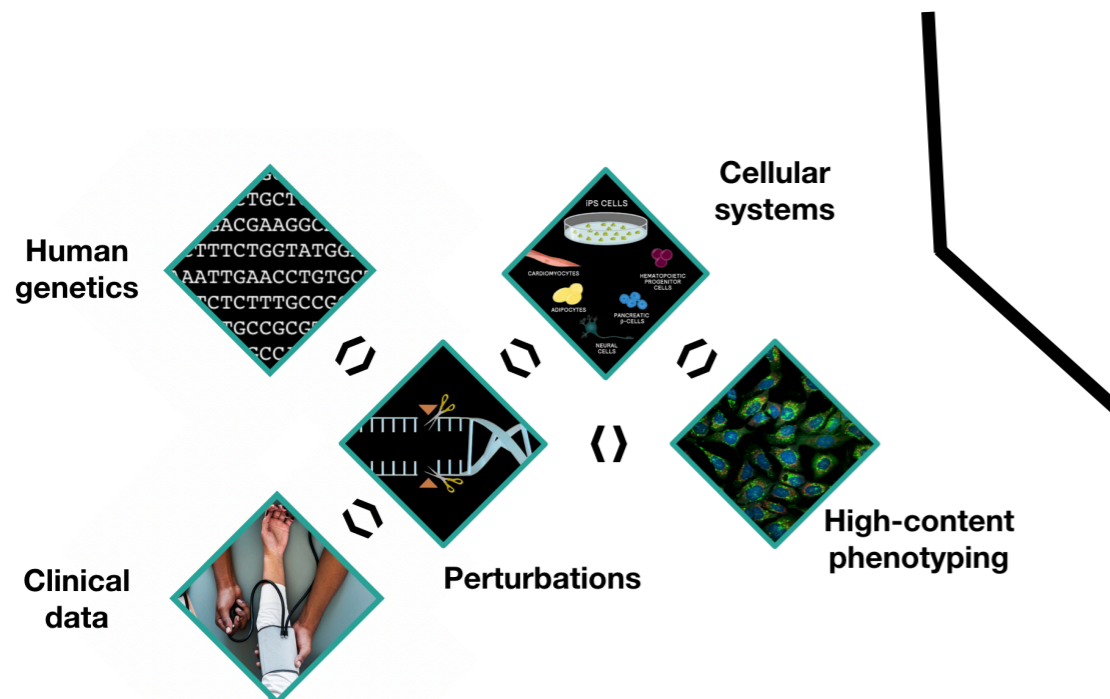
High-content microscopy



insitro is founded on data-driven technologies

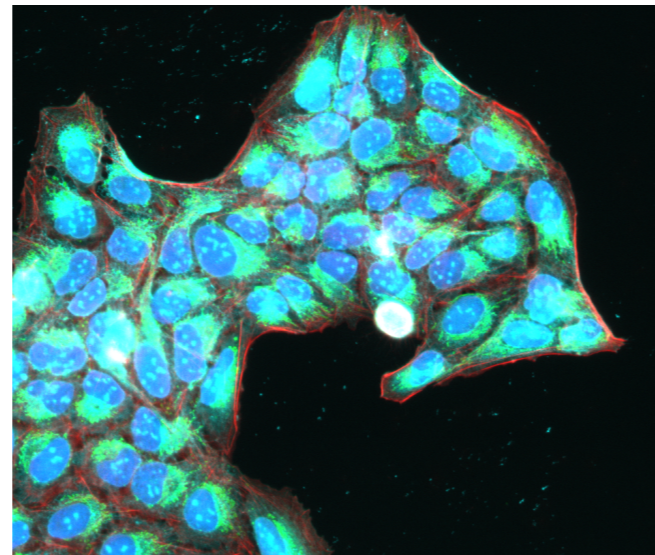


High-content microscopy



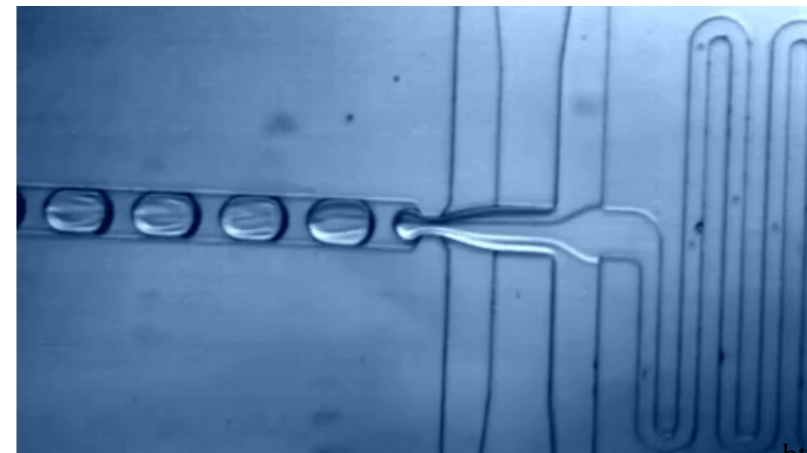
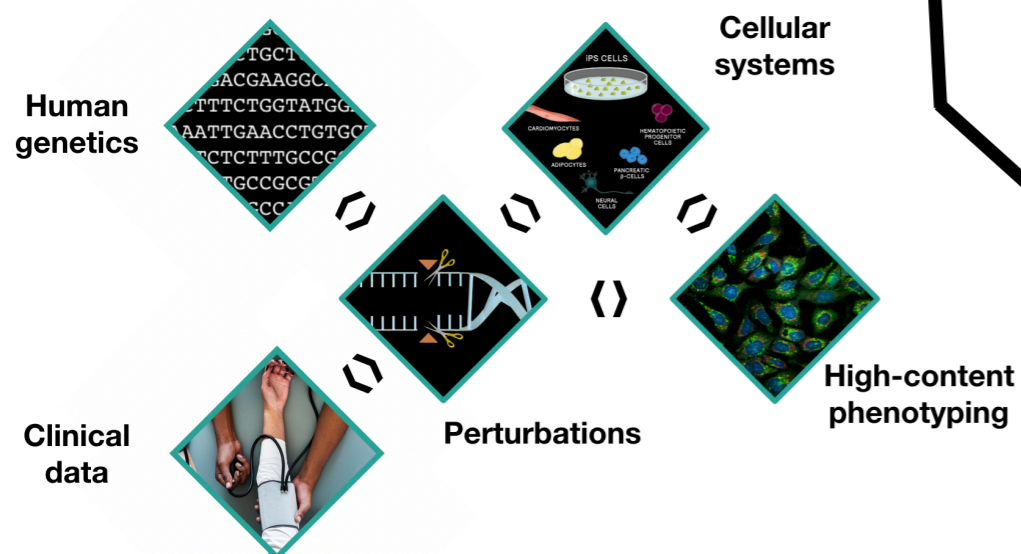
Single-cell RNA sequencing

insitro is founded on data-driven technologies



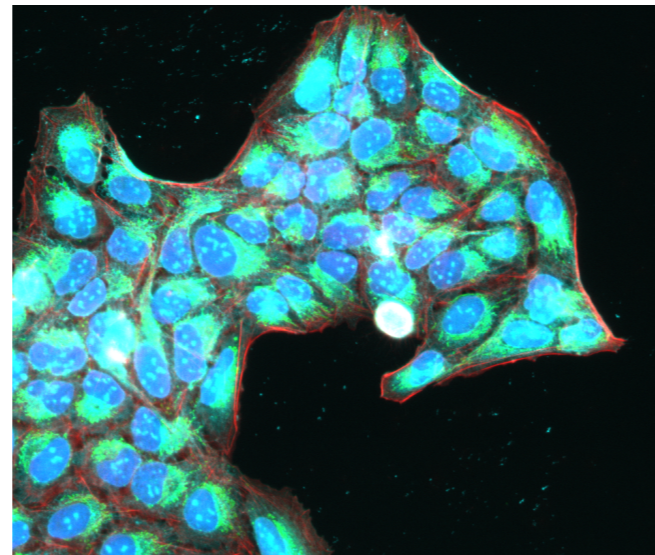
High-content microscopy

Machine Learning



Single-cell RNA sequencing

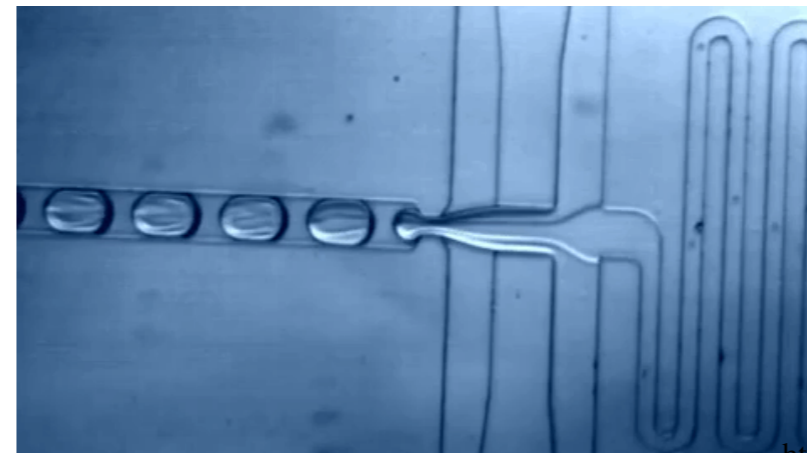
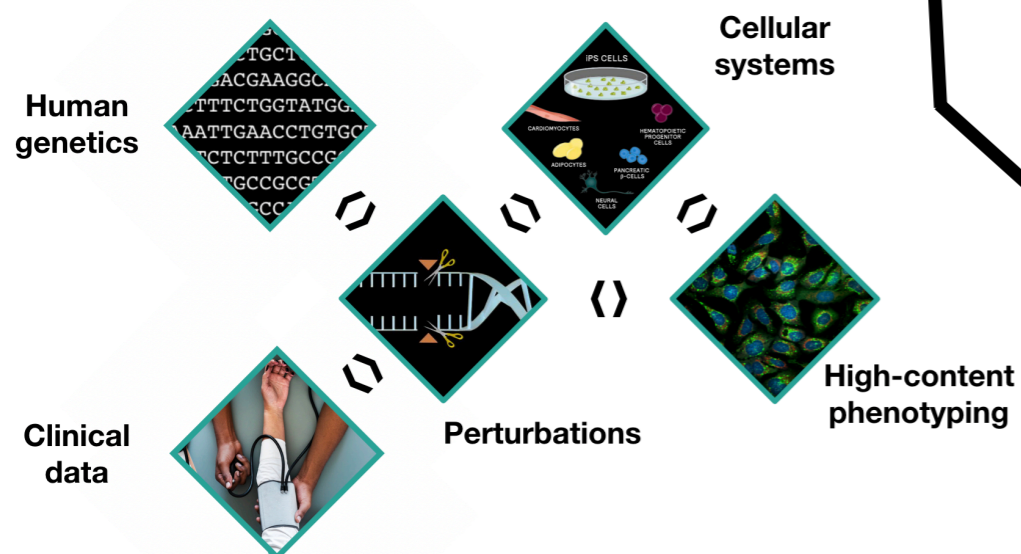
insitro is founded on data-driven technologies



High-content microscopy

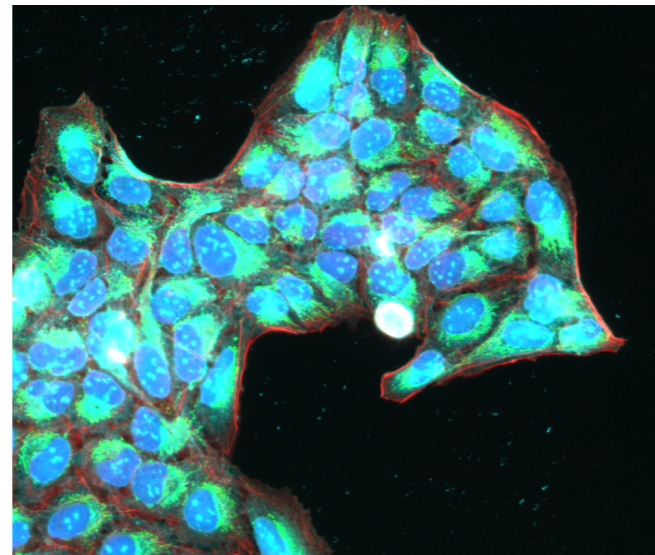
Machine Learning

Characterize phenotypes



Single-cell RNA sequencing

insitro is founded on data-driven technologies

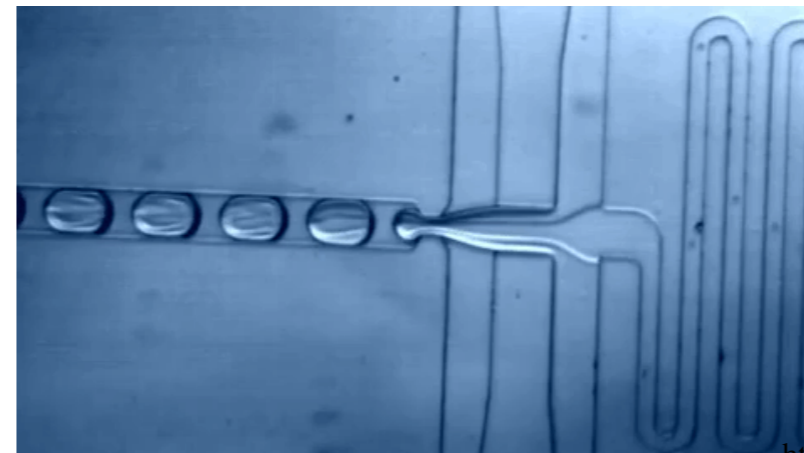
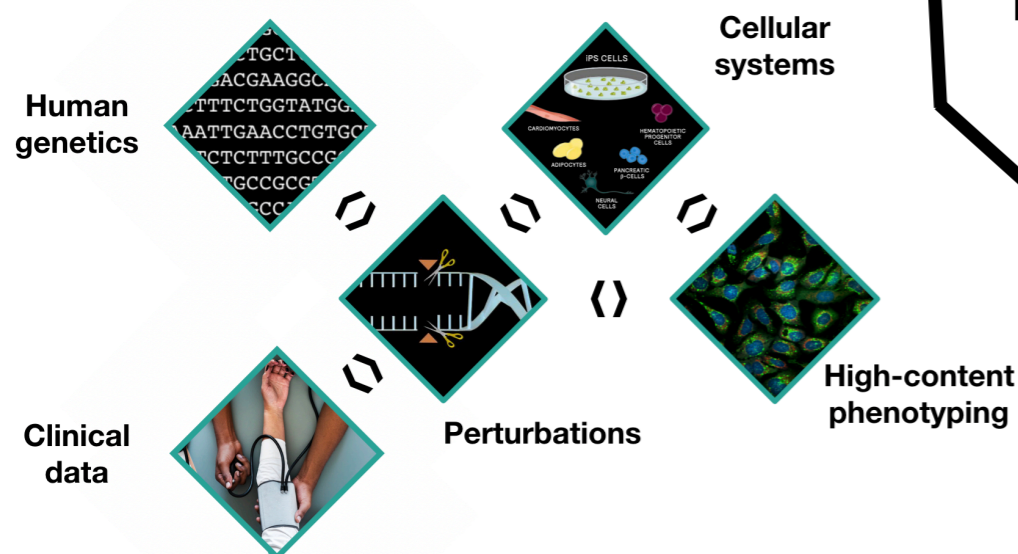


High-content microscopy

Machine Learning

Characterize phenotypes

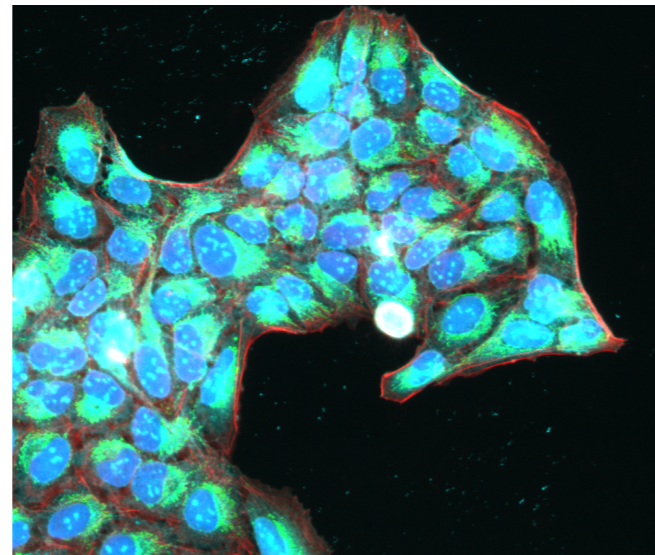
Understand relationships between data modalities



Single-cell RNA sequencing

insitro is founded on data-driven technologies

One experiment can test thousands of biological hypotheses

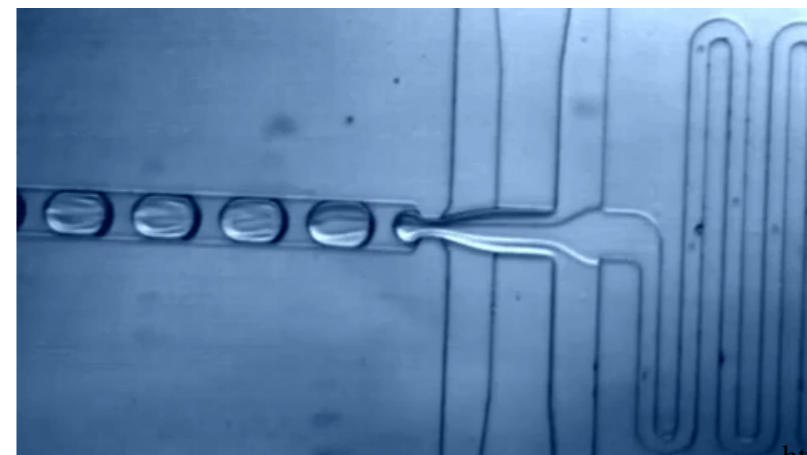
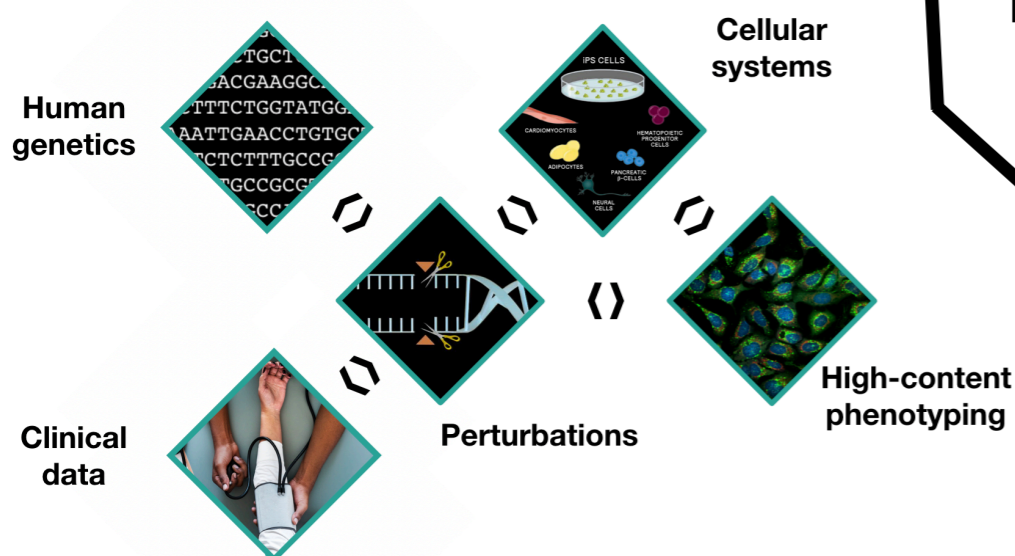


High-content microscopy

Machine Learning

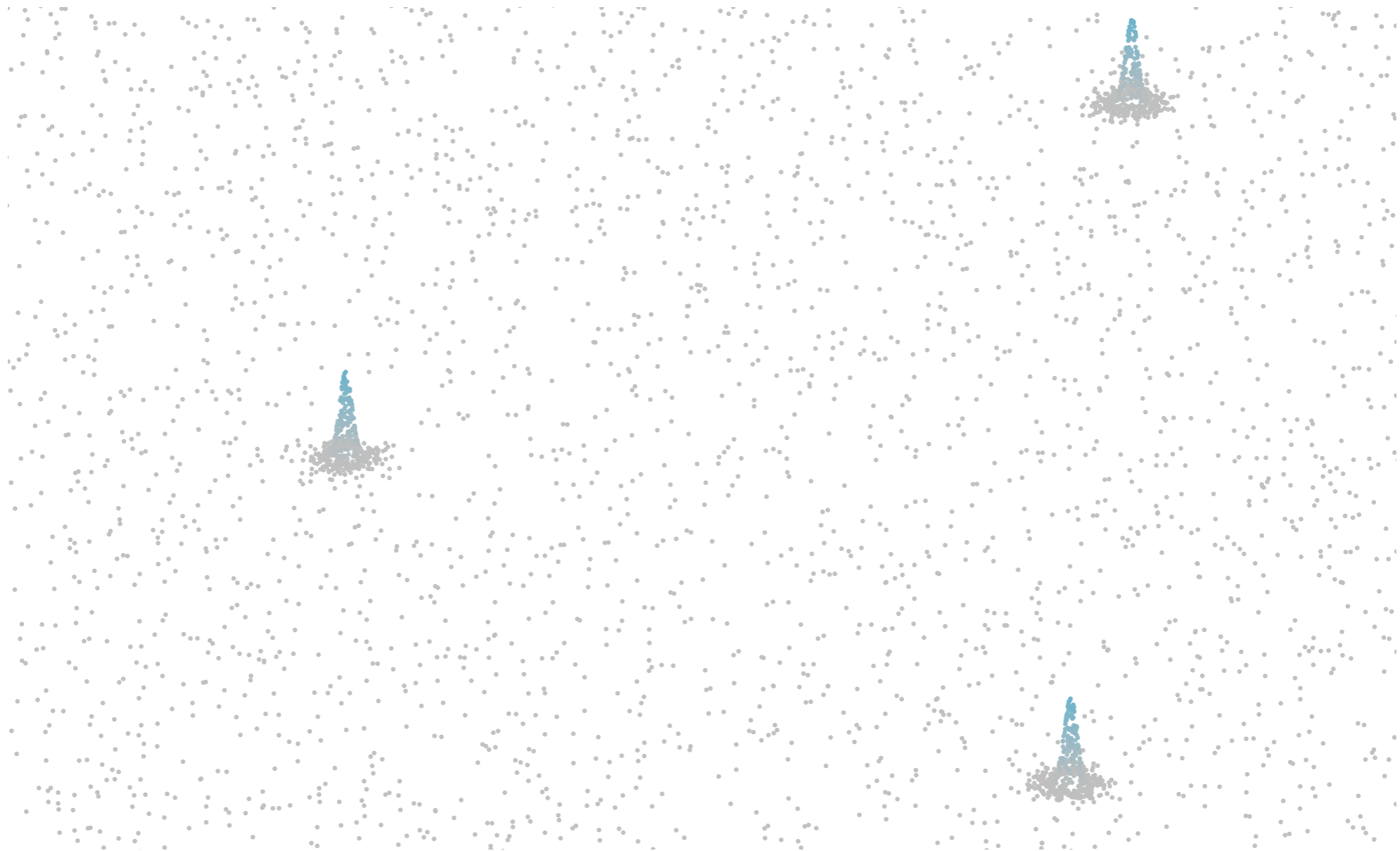
Characterize phenotypes

Understand relationships between data modalities

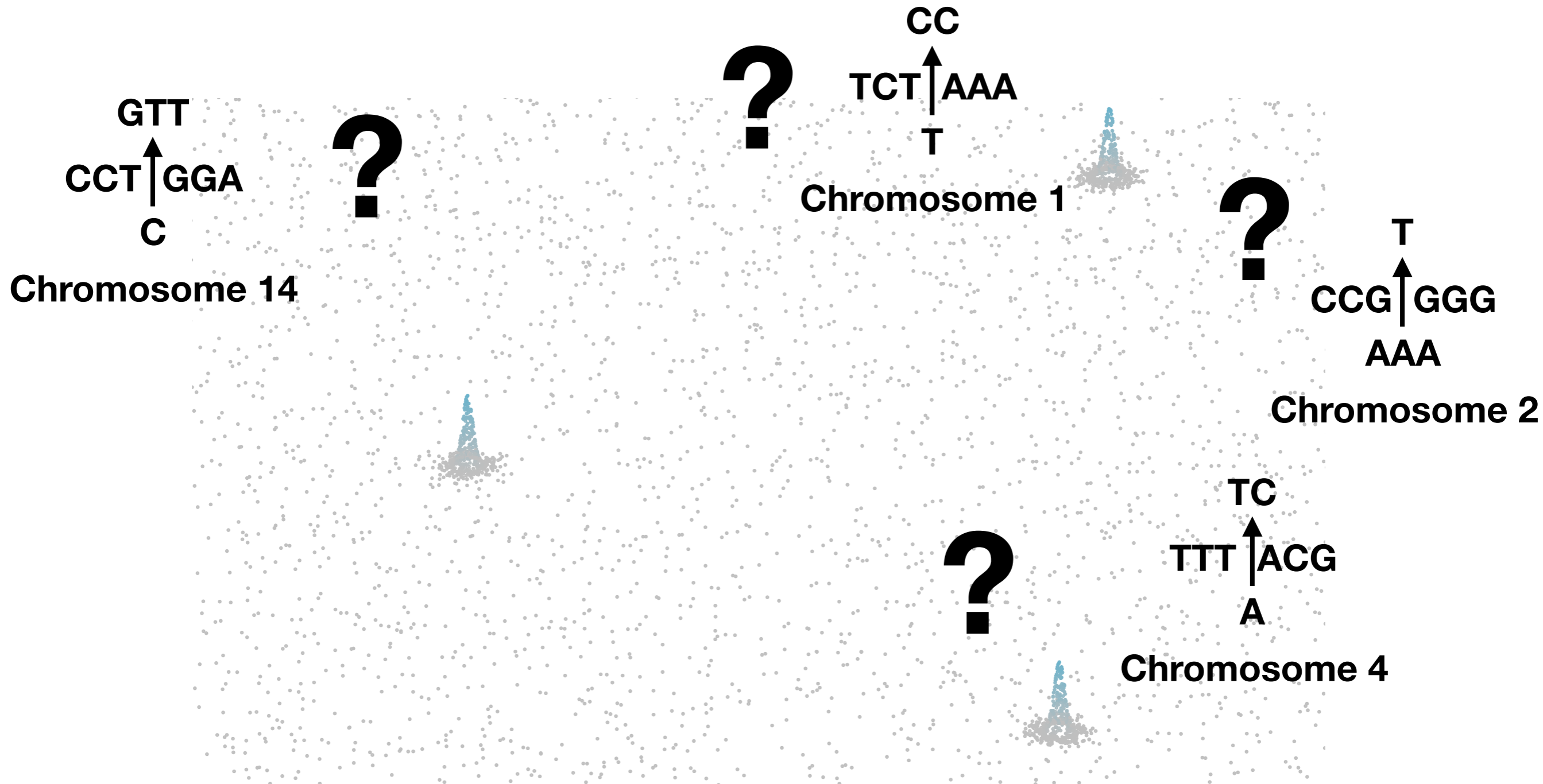


Single-cell RNA sequencing

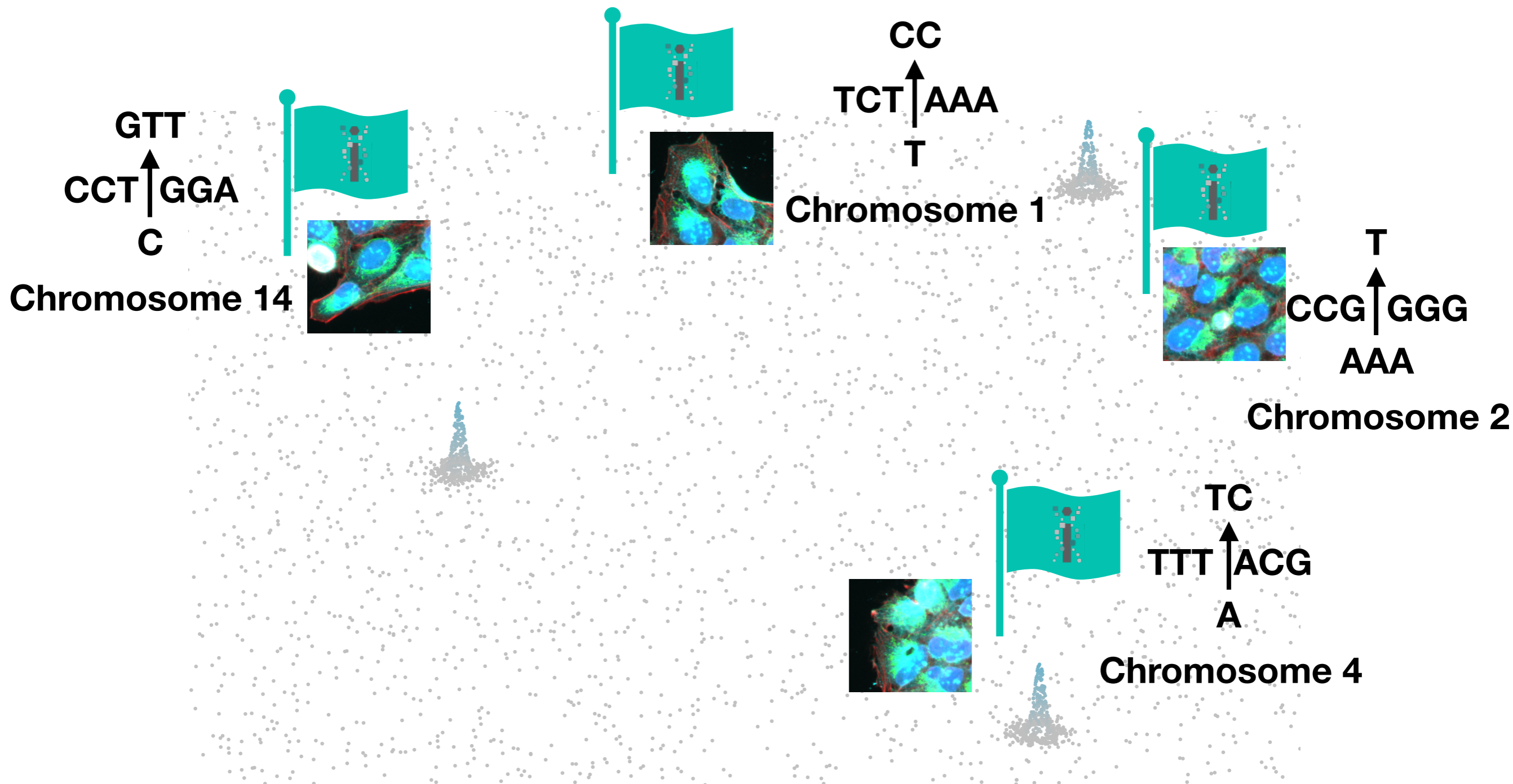
The insitro **bio-data factory** will elucidate the **landscape of disease**



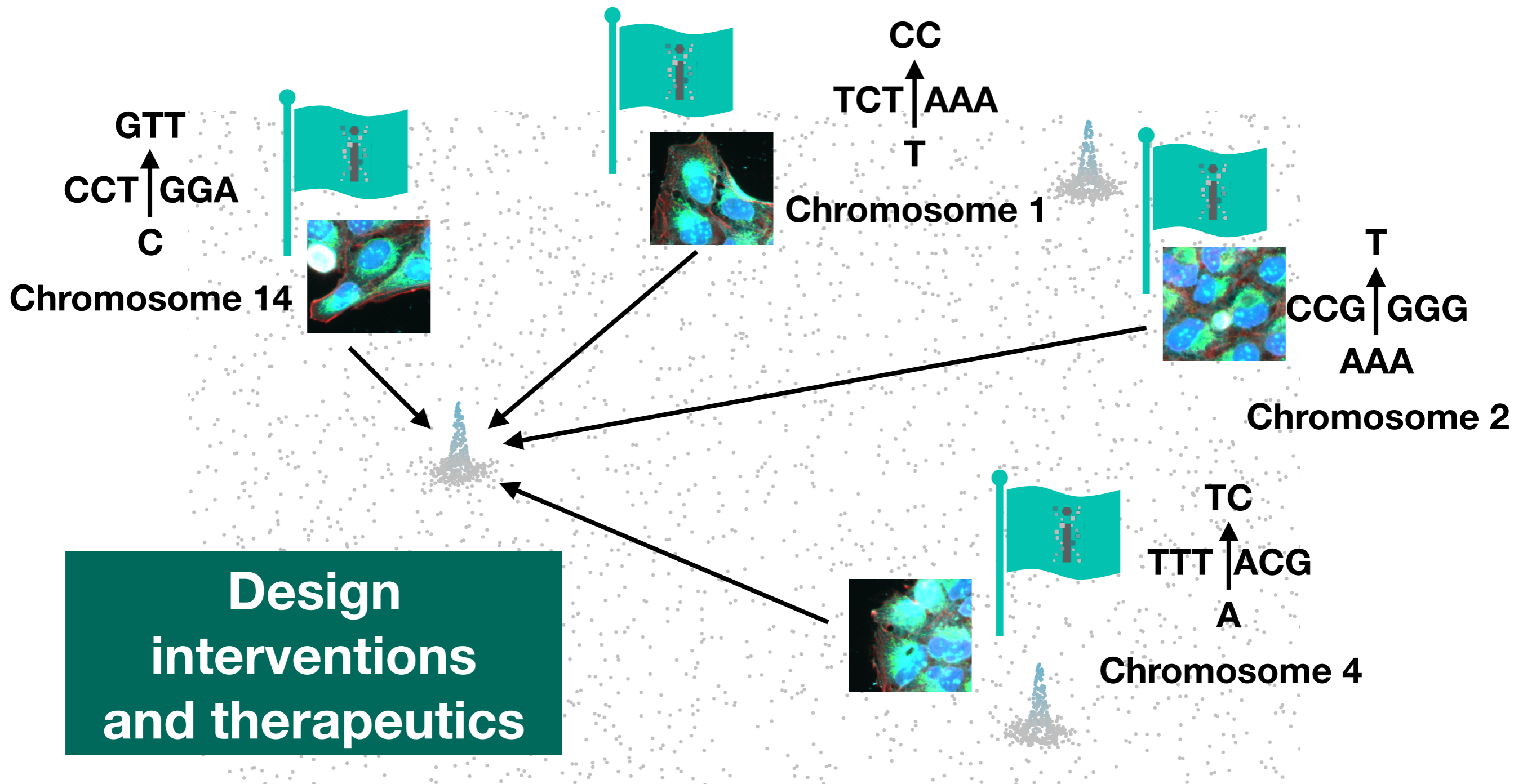
The insitro **bio-data factory** will elucidate the **landscape of disease**



The insitro **bio-data factory** will elucidate the **landscape of disease**



The insitro **bio-data factory** will elucidate the **landscape of disease**



The future is

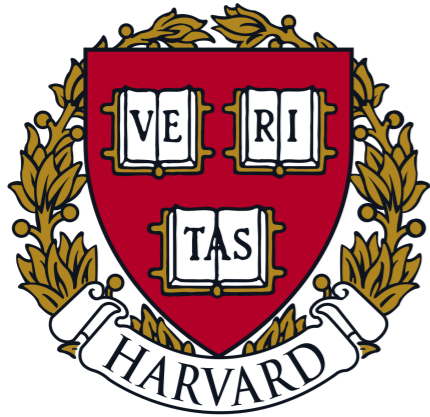
biology + computation

Thank you!

Thank you!



Thank you!



+



Debora Marks
John Ingraham
Chris Sander
Andrew Kruse
Aashish Manglik
Conor McMahon
June Shin
Aaron Kollasch

Anna Green
Thomas Hopf
Charlotta Scharfe
Benni Schubert
Eli Weinstein
Kelly Brock
Rohan Maddamsetti
David Ding
Hailey Cambra
Agnes Toth-Petroczy
Perry Palmedo
Frank Poelwijk
Nick Gauthier
Jennie Epp



Sam Deutsch



Daphne Koller

