

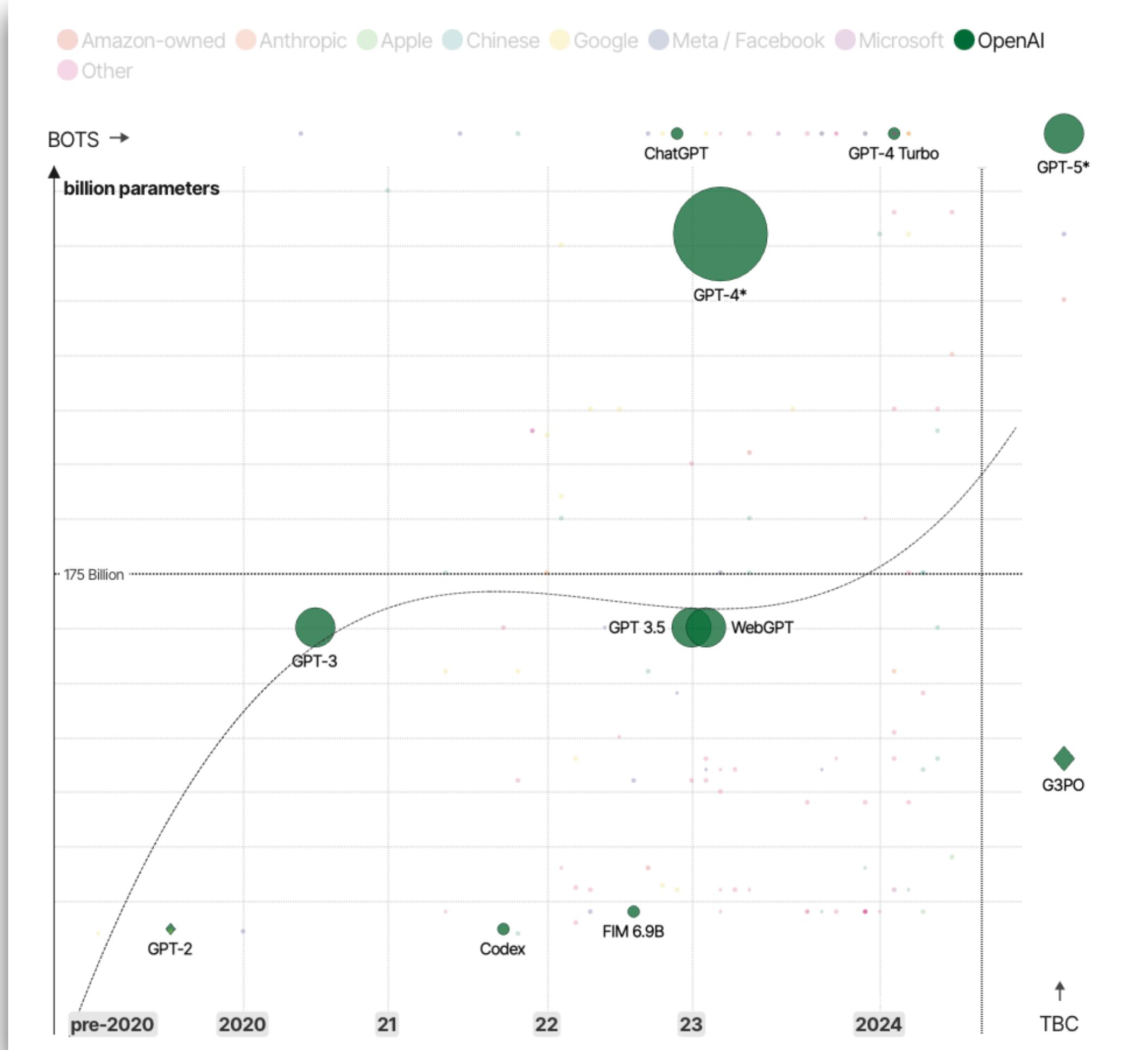
# SHAP zero explains biological sequence models with near-zero marginal cost for future queries

Darin Tsui

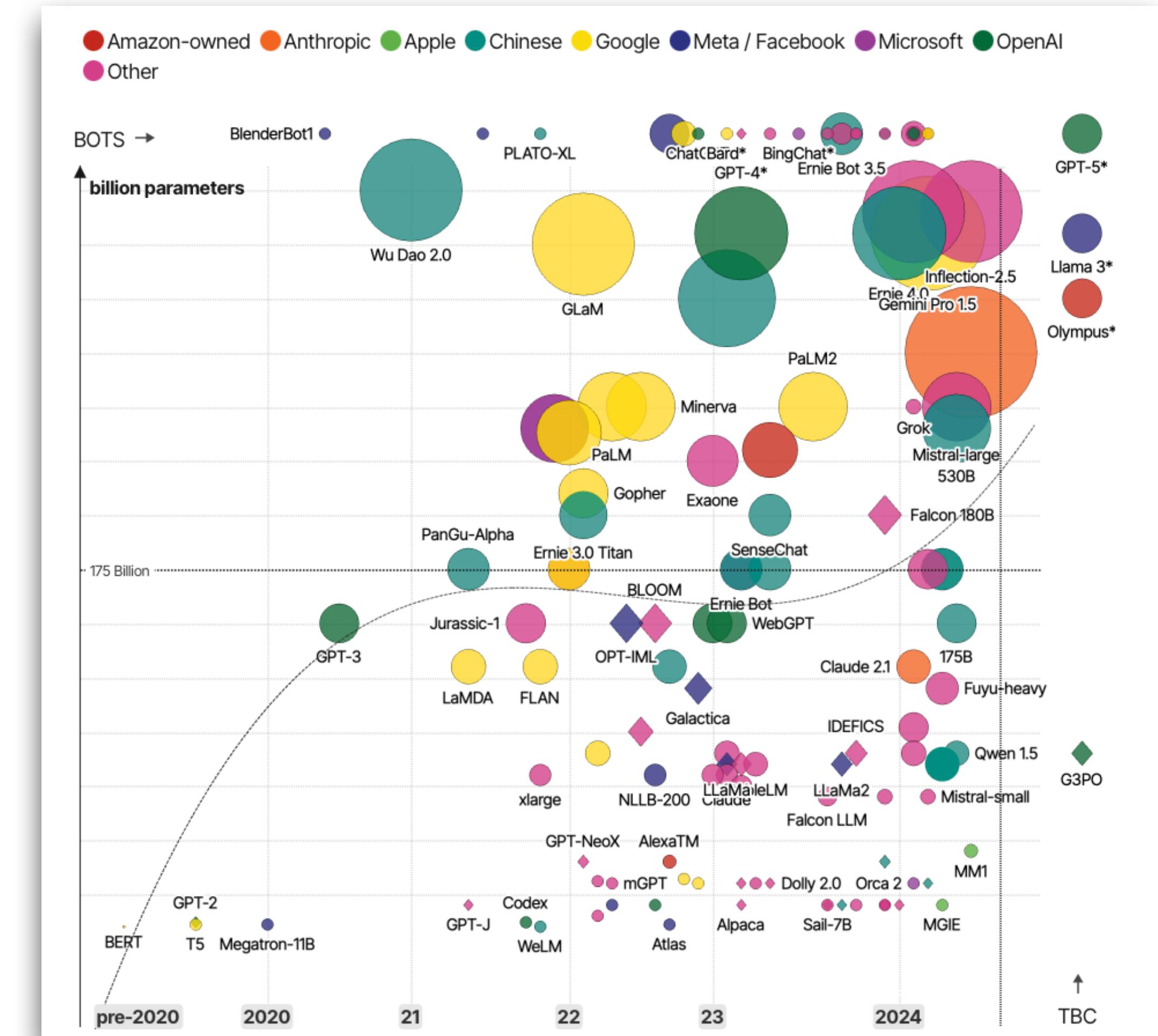
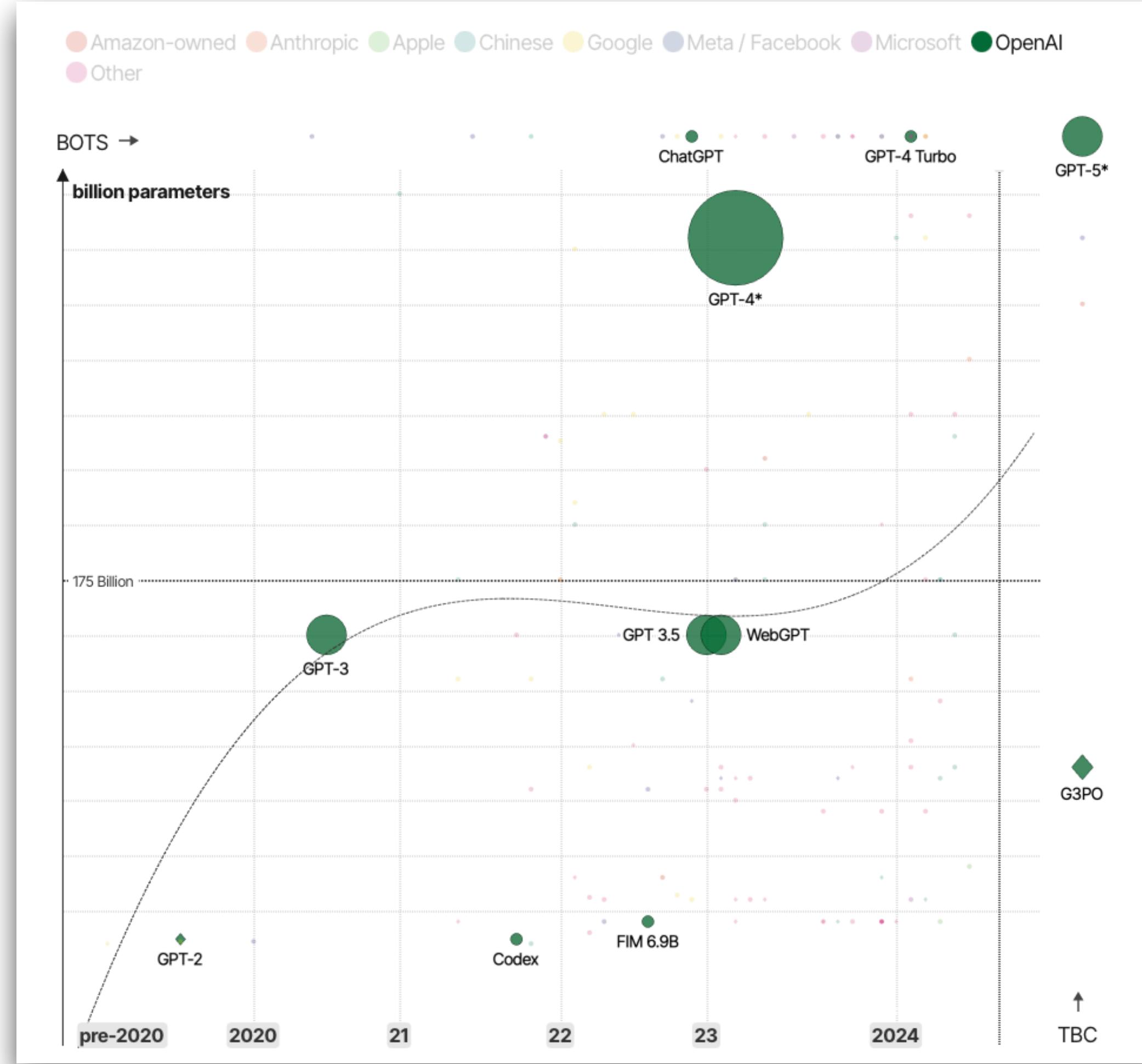


Georgia Tech College of Engineering  
**School of Electrical  
and Computer Engineering**

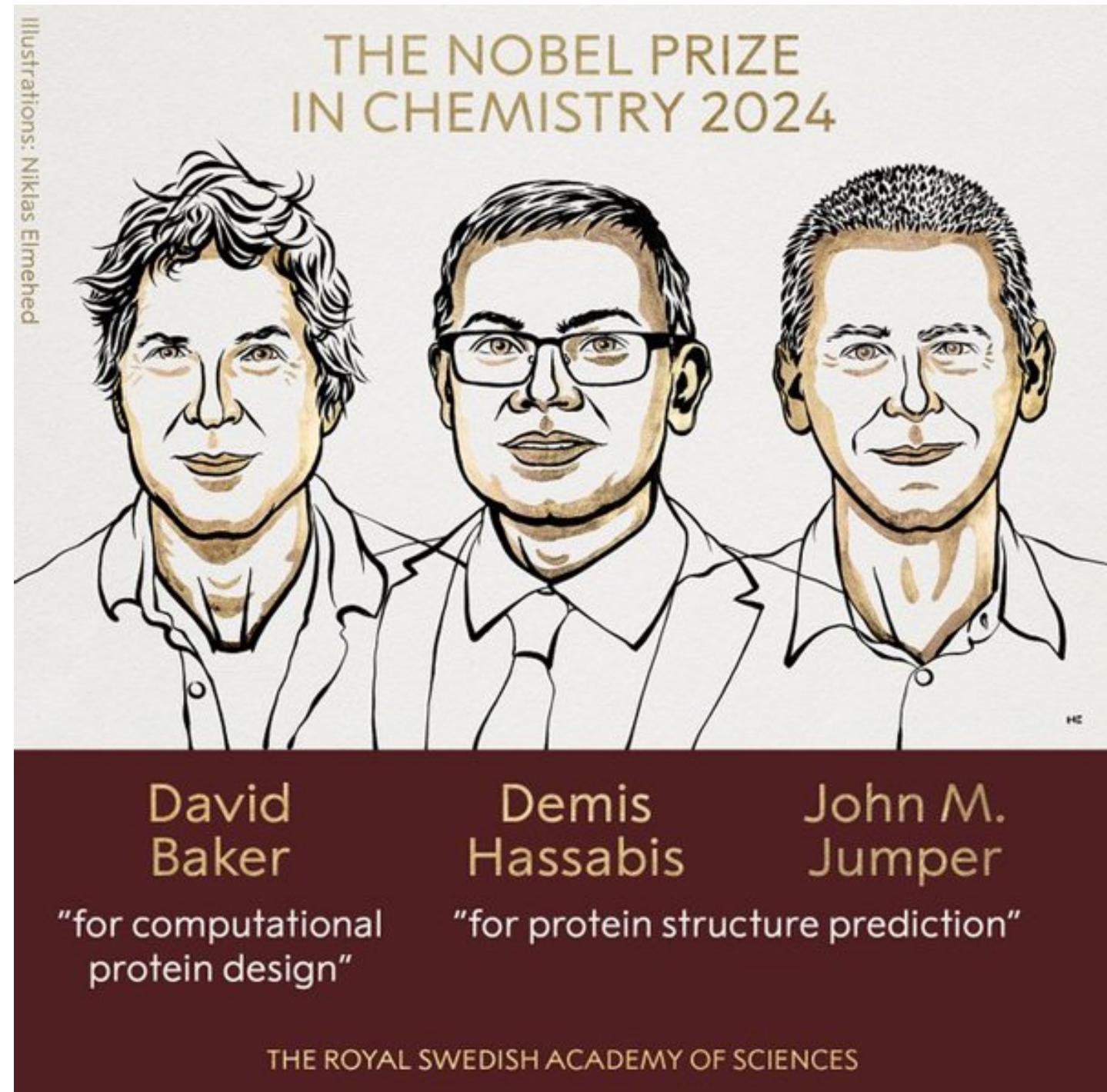
# the rise and rise of AI



# the rise and rise of AI



# AI breakthroughs in science



# AI breakthroughs in science



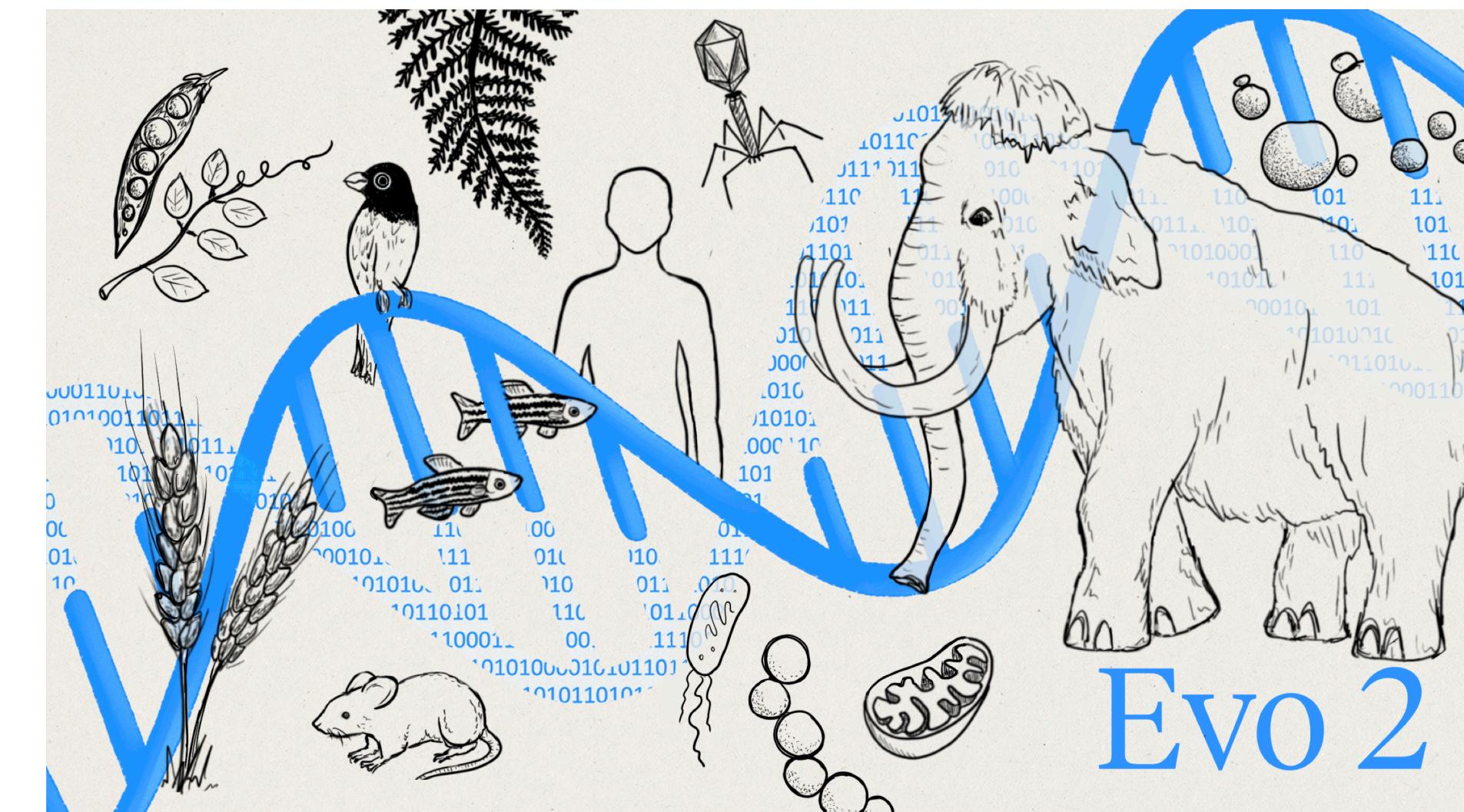
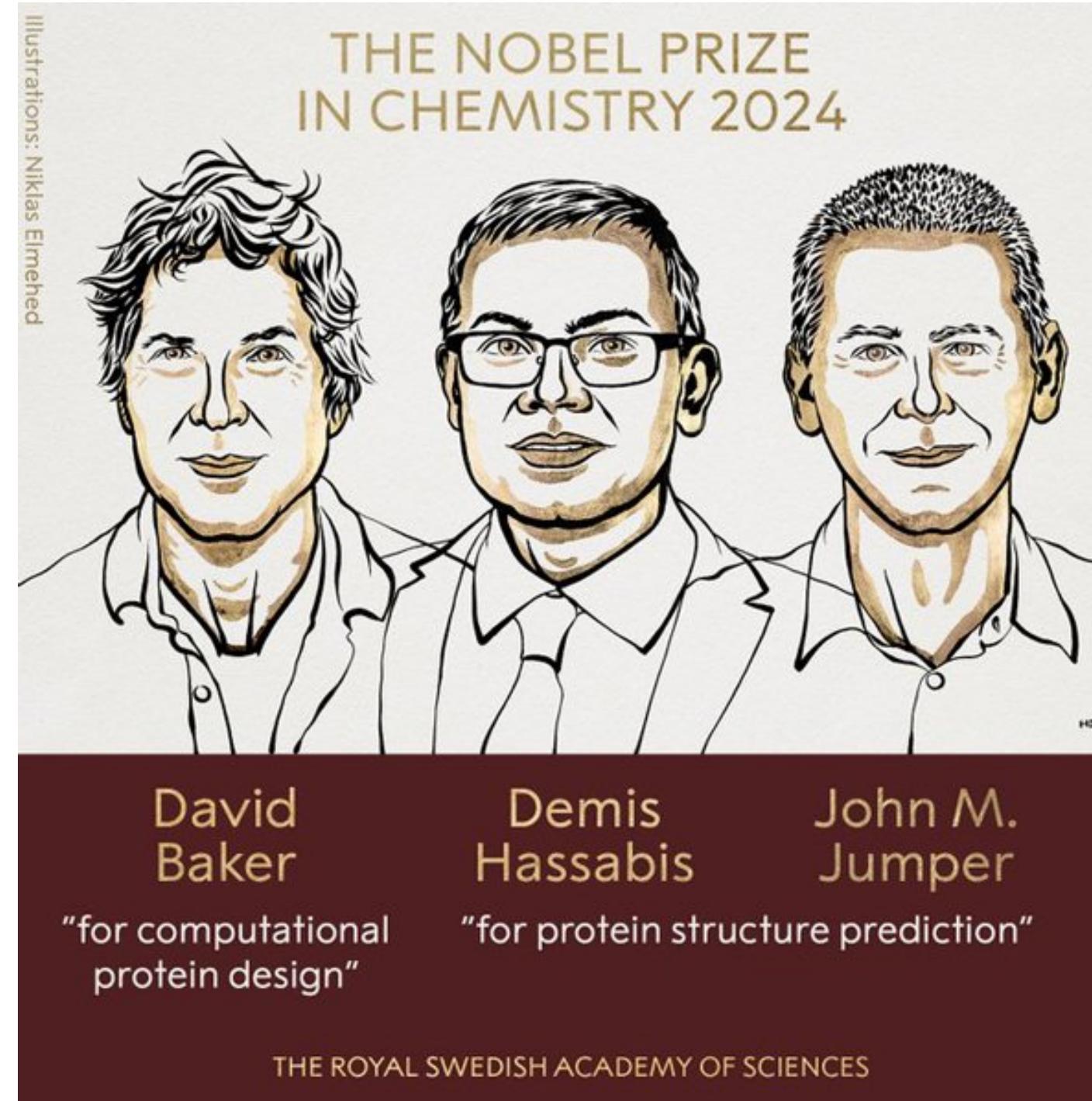
JUNE 25, 2024//ANNOUNCEMENT

## ESM3: Simulating 500 million years of evolution with a language model

[Preview our paper ➔](#)

Trained on **2.78 billion** natural proteins  
**96 billion** parameters

# AI breakthroughs in science



Trained on **9.3 trillion nucleotides**  
**40 billion parameters**

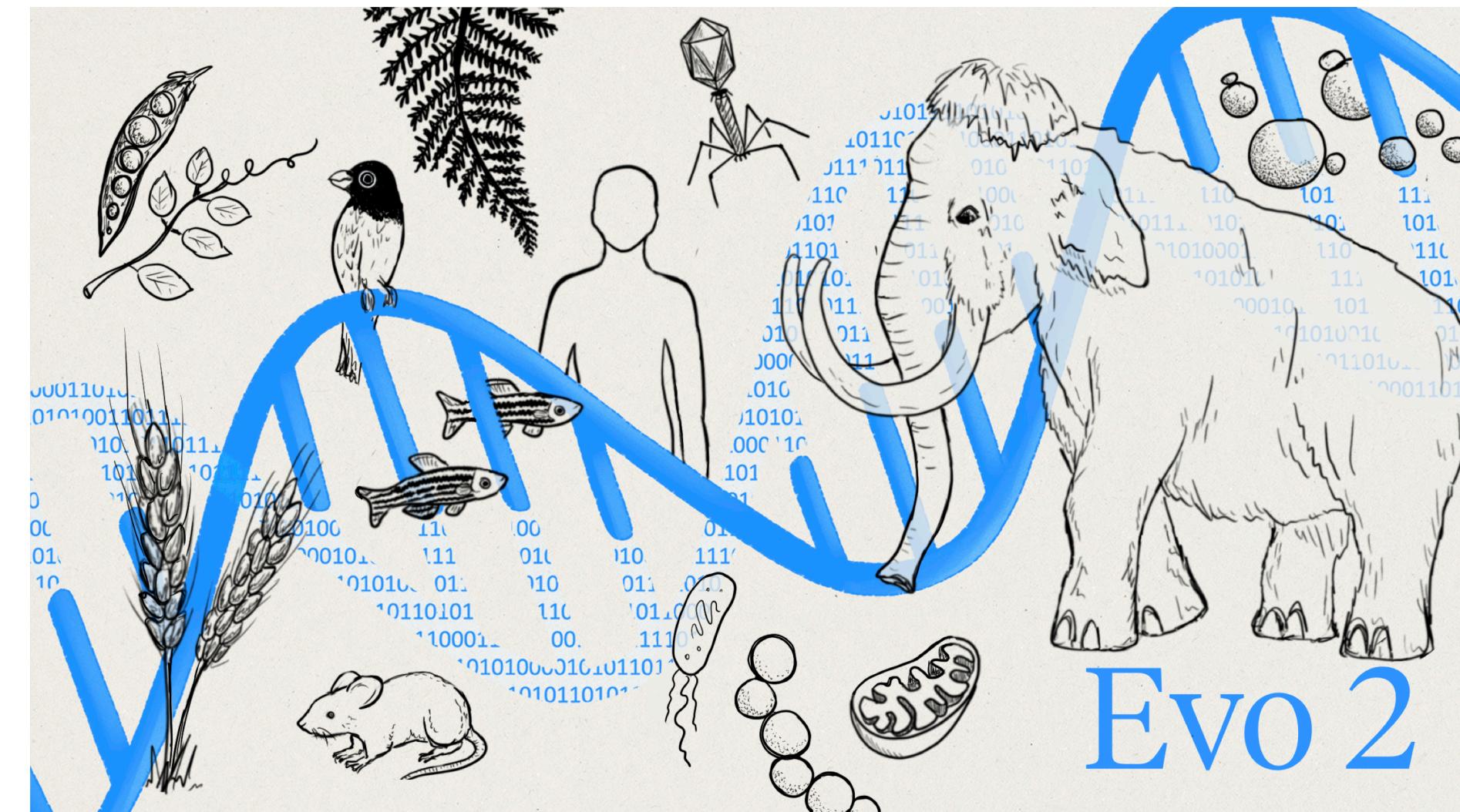
JUNE 25, 2024//ANNOUNCEMENT

**ESM3: Simulating 500 million years  
of evolution with a language model**

[Preview our paper ➔](#)

Trained on **2.78 billion natural proteins**  
**96 billion parameters**

# AI breakthroughs in science



Trained on **9.3 trillion nucleotides**  
**40 billion parameters**

JUNE 25, 2024//ANNOUNCEMENT

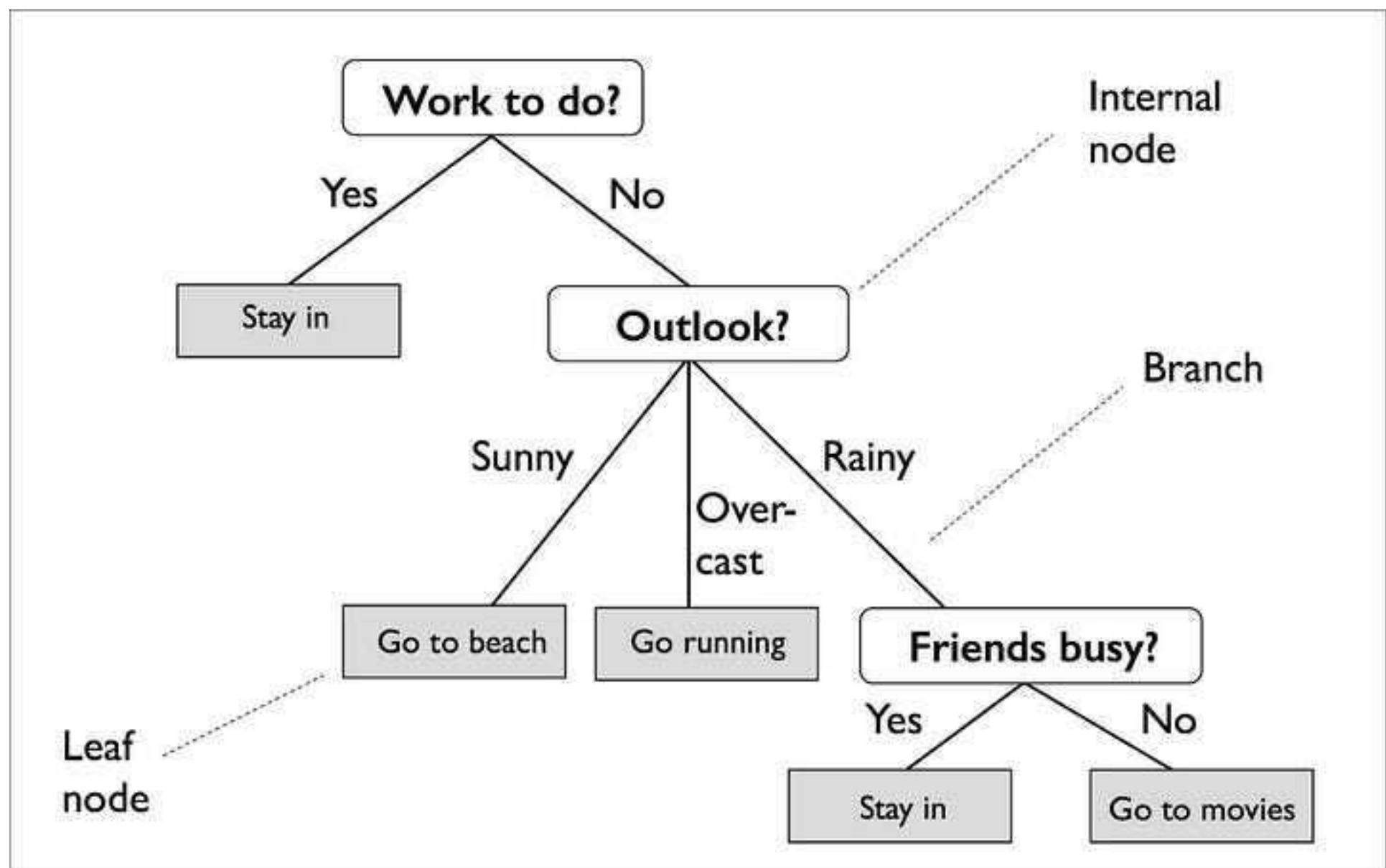
**ESM3: Simulating 500 million years  
of evolution with a language model**

[Preview our paper →](#)

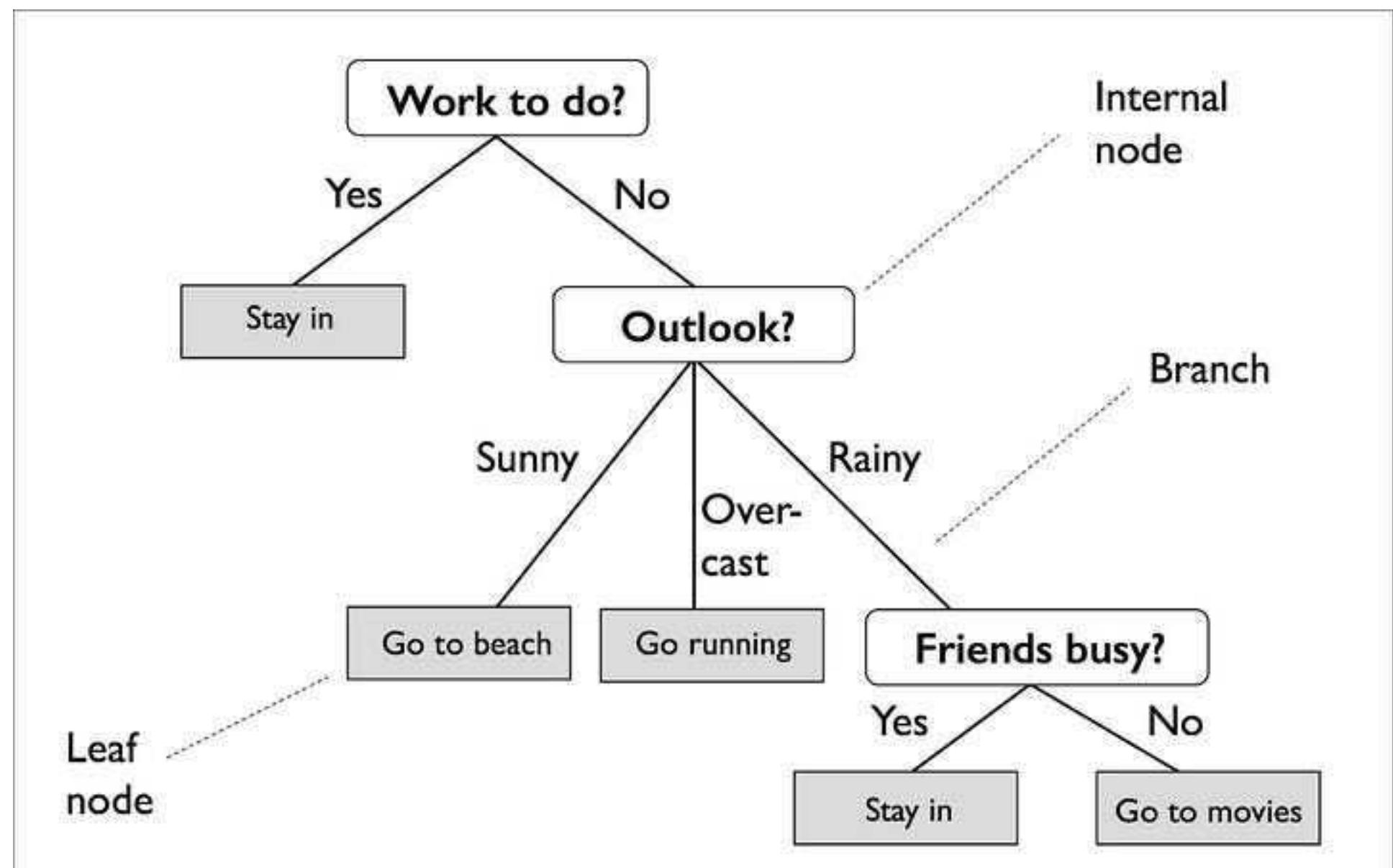
Trained on **2.78 billion natural proteins**  
**96 billion parameters**

But does our power to  
**explain** them also **scale**?

# interpretable machine learning models

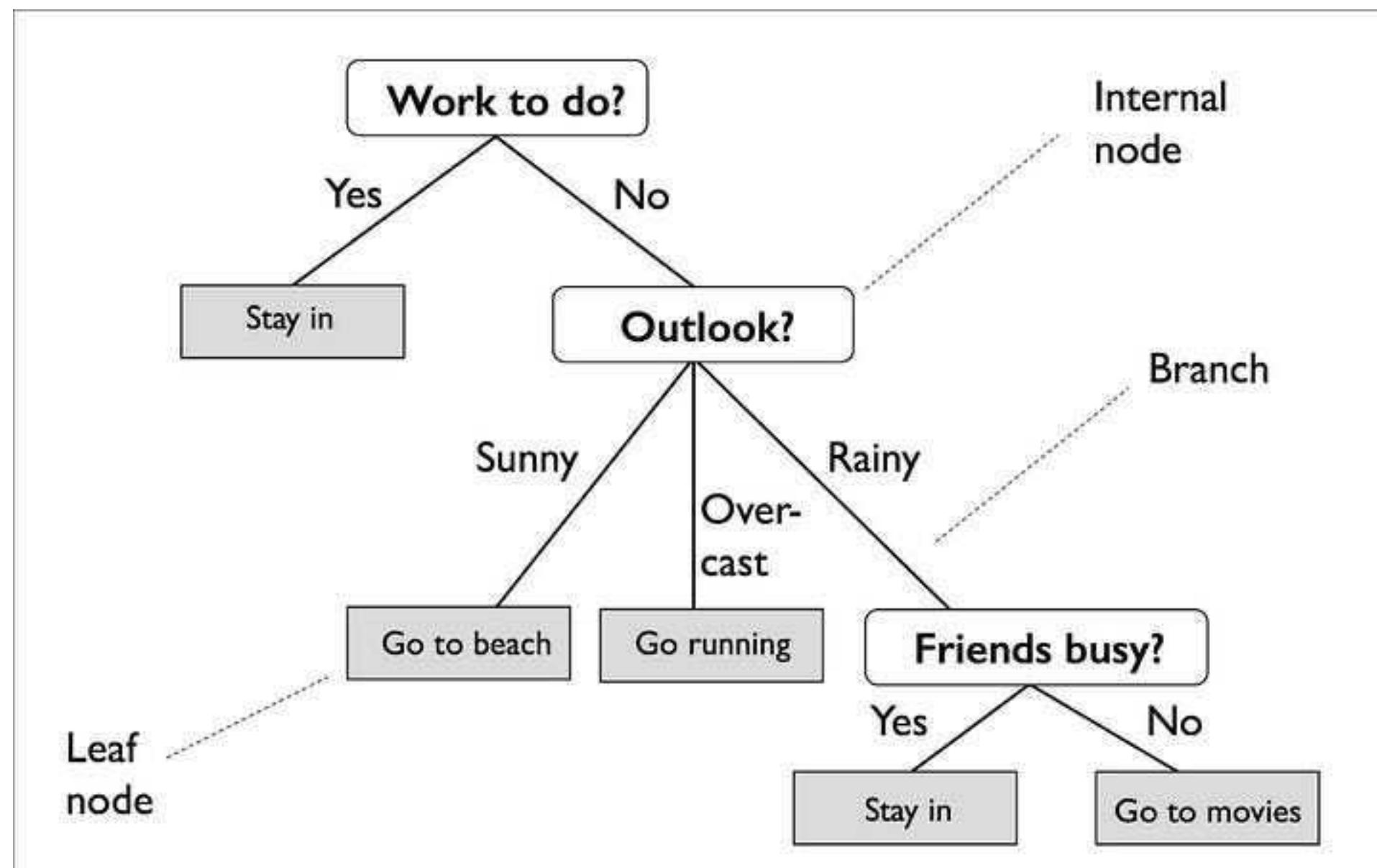


# interpretable machine learning models

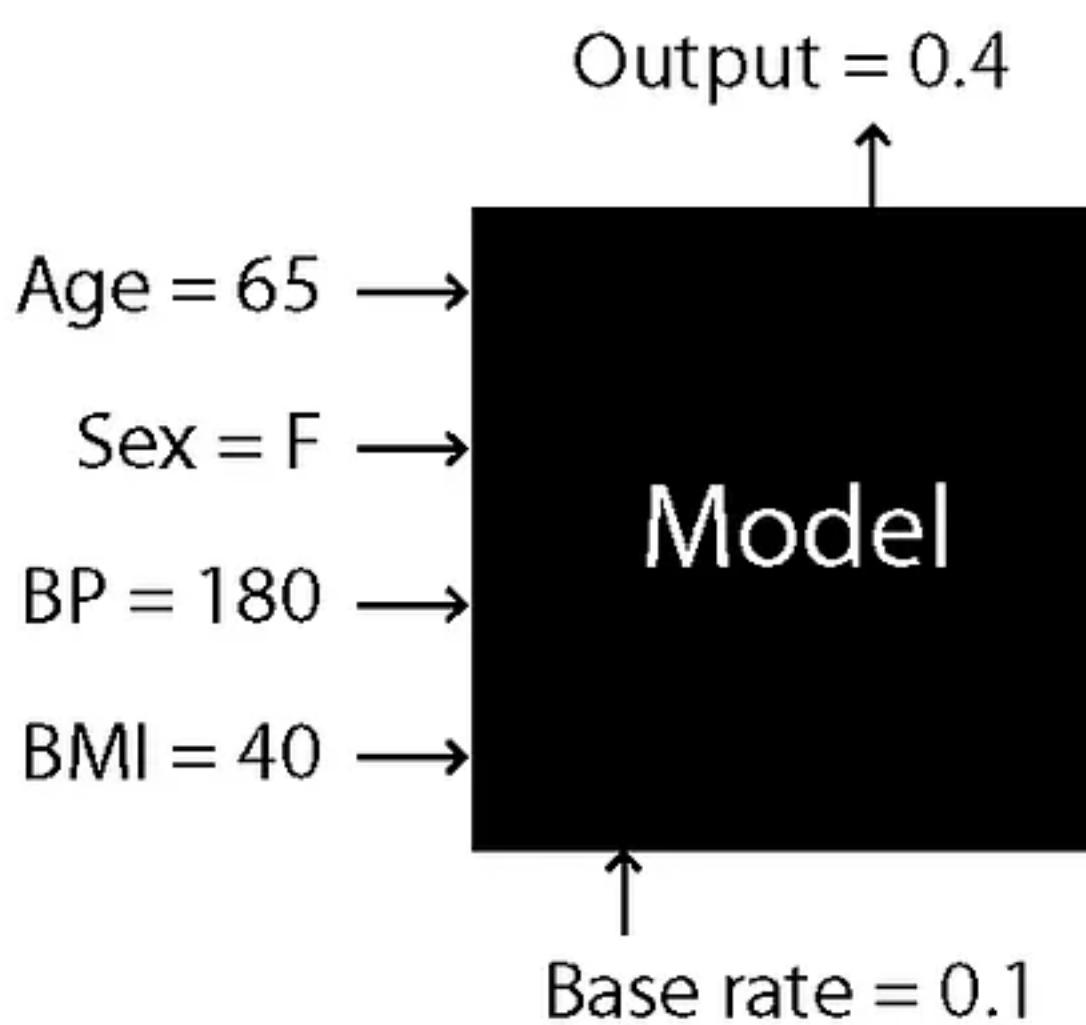


$$y = 3x_1 - x_2 + 2x_3$$

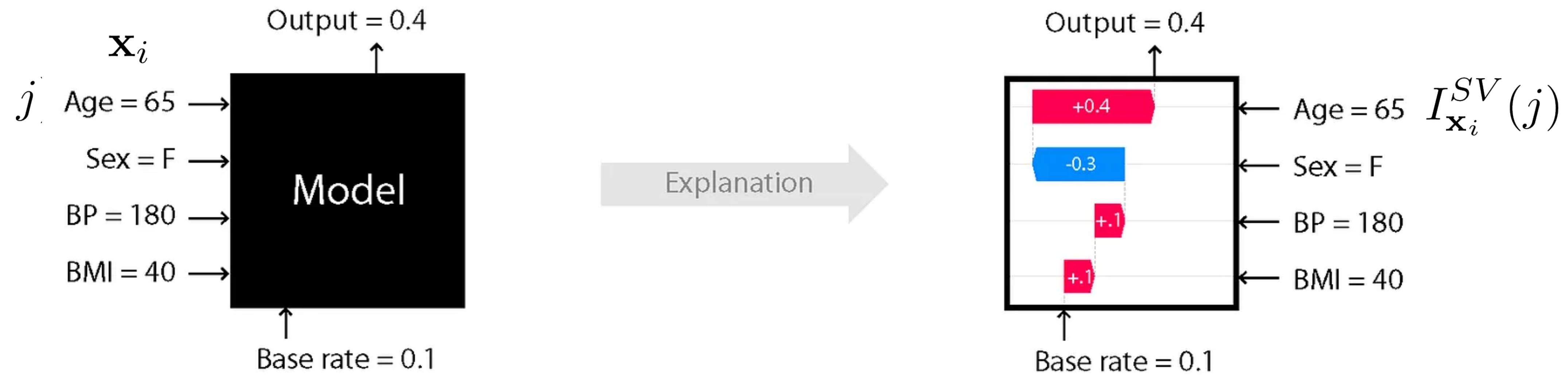
# interpretable machine learning models



$$y = 3x_1 - x_2 + 2x_3$$

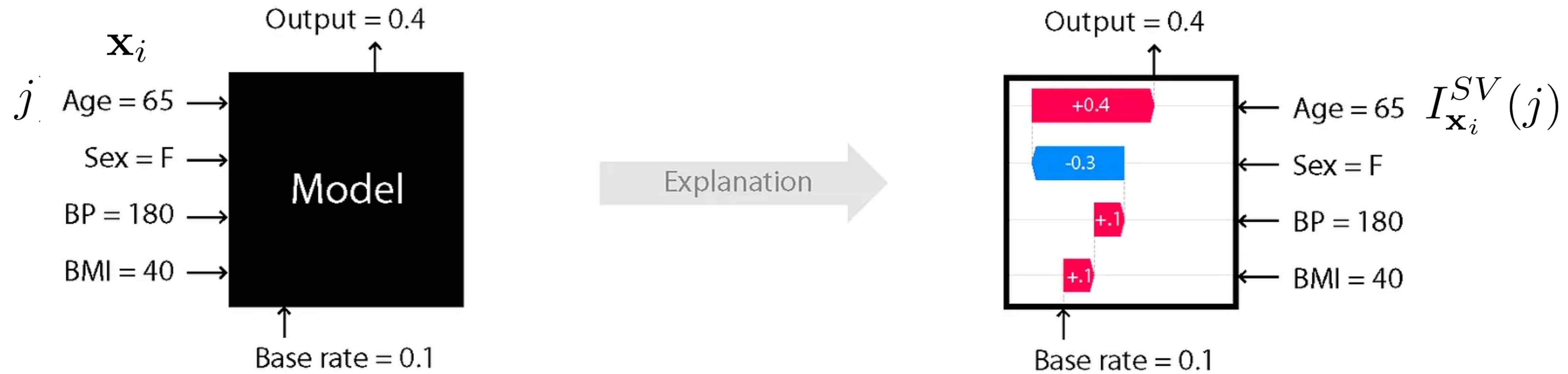


# Shapley values



$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

# Shapley values



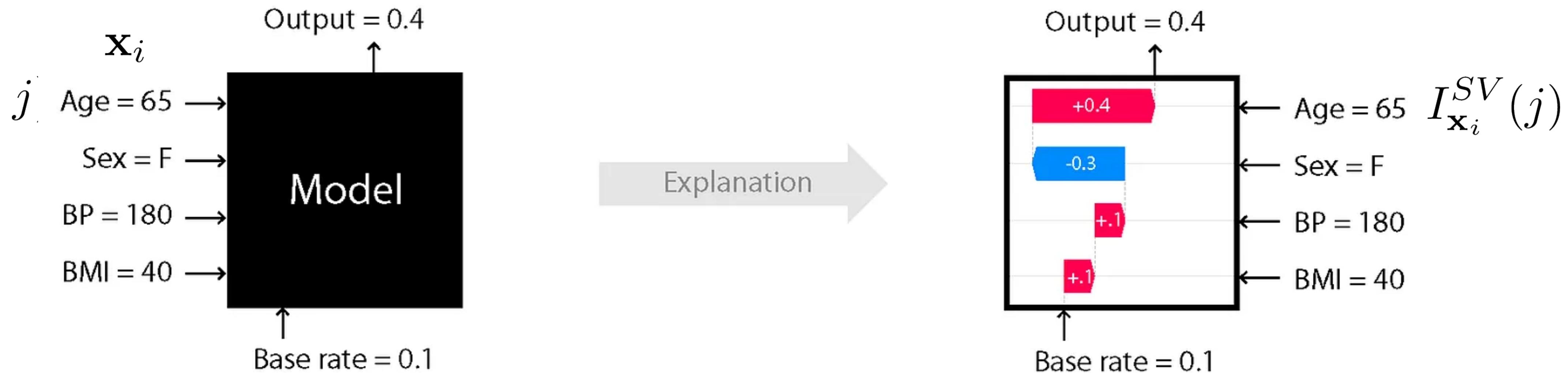
$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

$$\begin{aligned} I_{\mathbf{x}_i}^{SV}(j = Age) &= \frac{1}{4} [v_{\{Age\}}(\mathbf{x}_i) - v_{\emptyset}(\mathbf{x}_i)] \\ &\quad + \frac{1}{12} [v_{\{Age, Sex\}}(\mathbf{x}_i) - v_{\{Sex\}}(\mathbf{x}_i)] \end{aligned}$$

...

$$+ \frac{1}{4} [v_{\{Age, Sex, BP, BMI\}}(\mathbf{x}_i) - v_{\{Sex, BP, BMI\}}(\mathbf{x}_i)]$$

# Shapley values



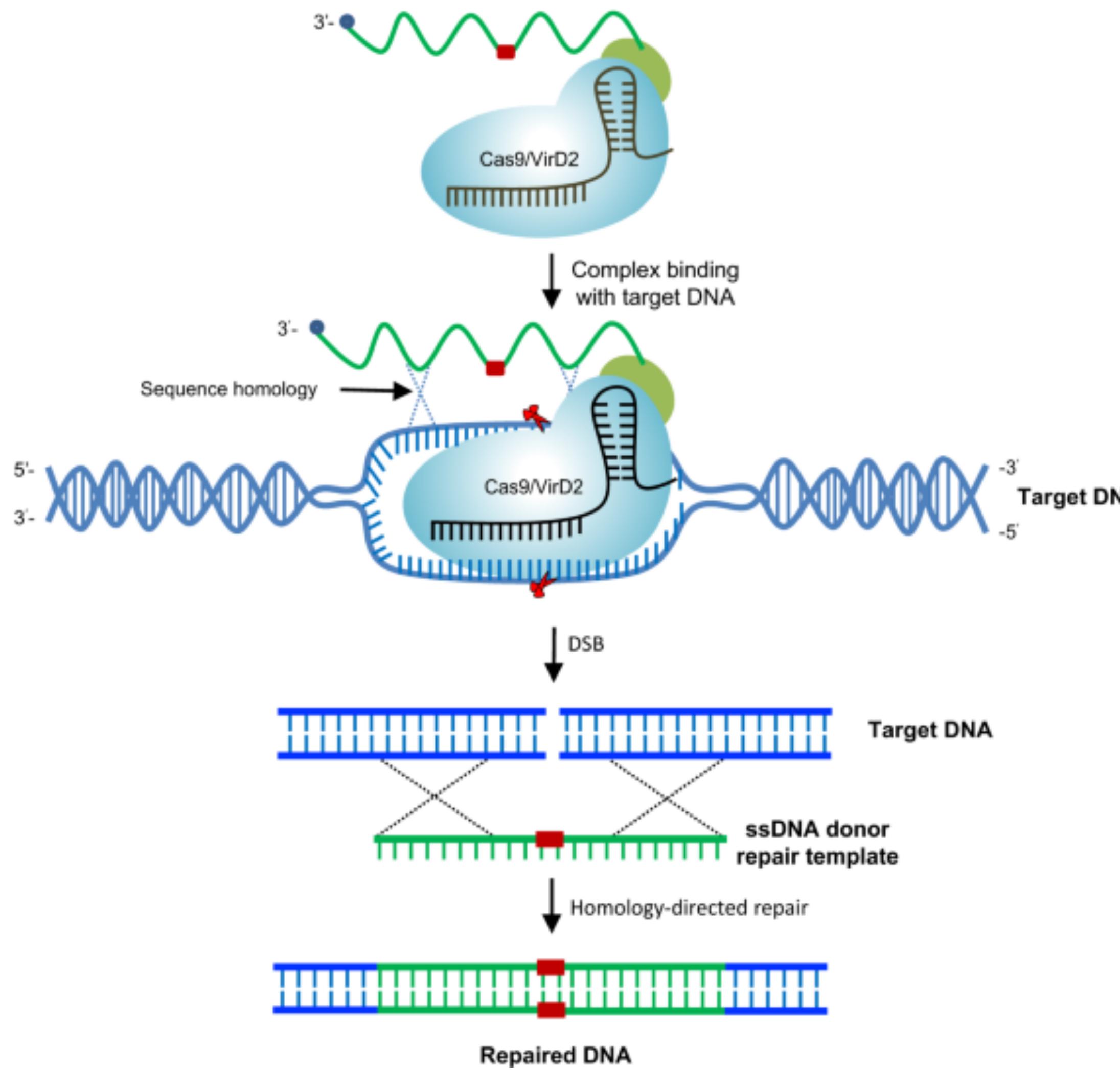
$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

$$I_{\mathbf{x}_i}^{SV}(j = Age) = \frac{1}{4} [v_{\{Age\}}(\mathbf{x}_i) - v_{\emptyset}(\mathbf{x}_i)]$$

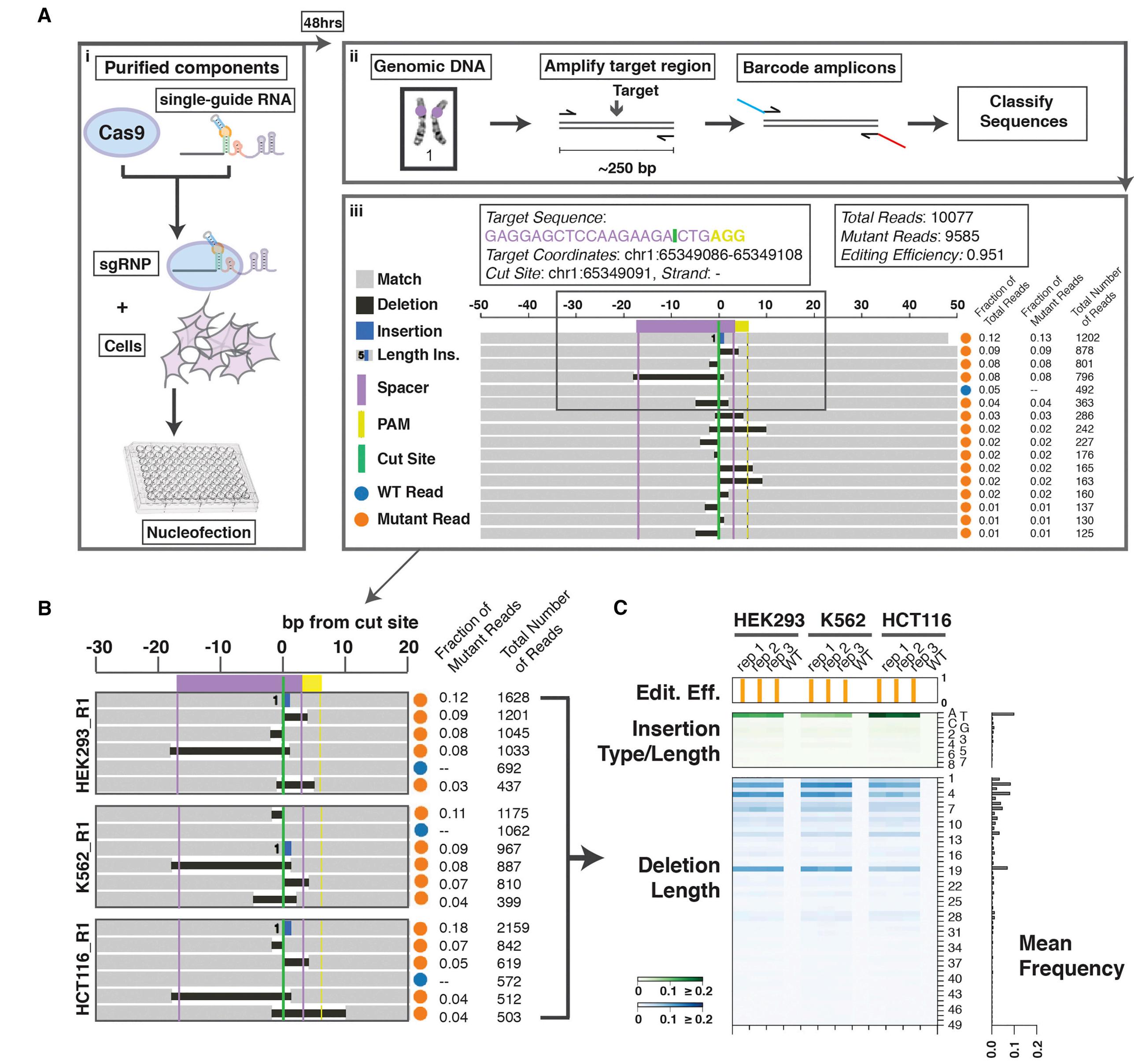
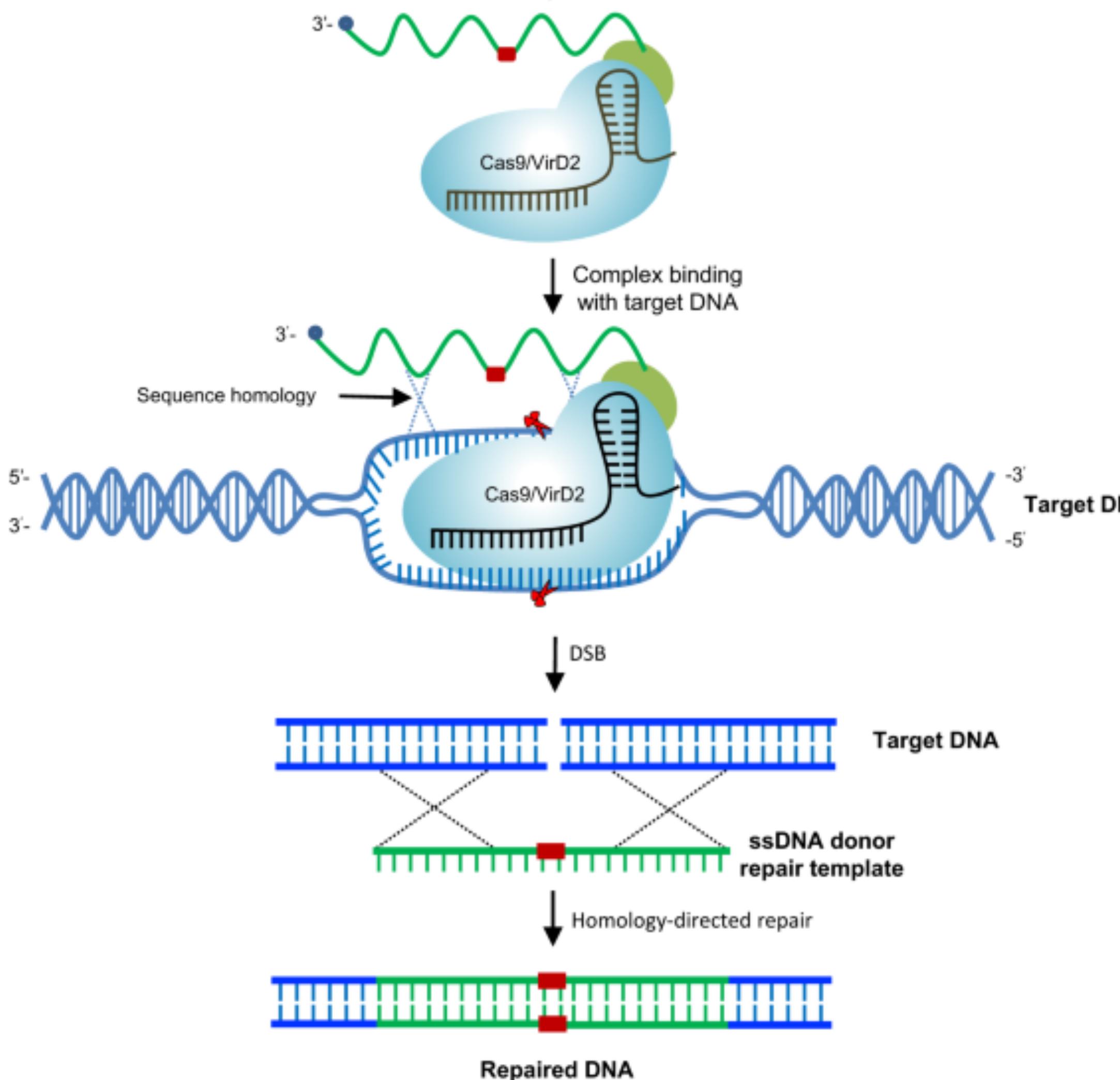
+  $\frac{1}{4}$  [But how about explaining  
many sequences?]

$$+ \frac{1}{4} [v_{\{Age, Sex, BP, BMI\}}(\mathbf{x}_i) - v_{\{Sex, BP, BMI\}}(\mathbf{x}_i)]$$

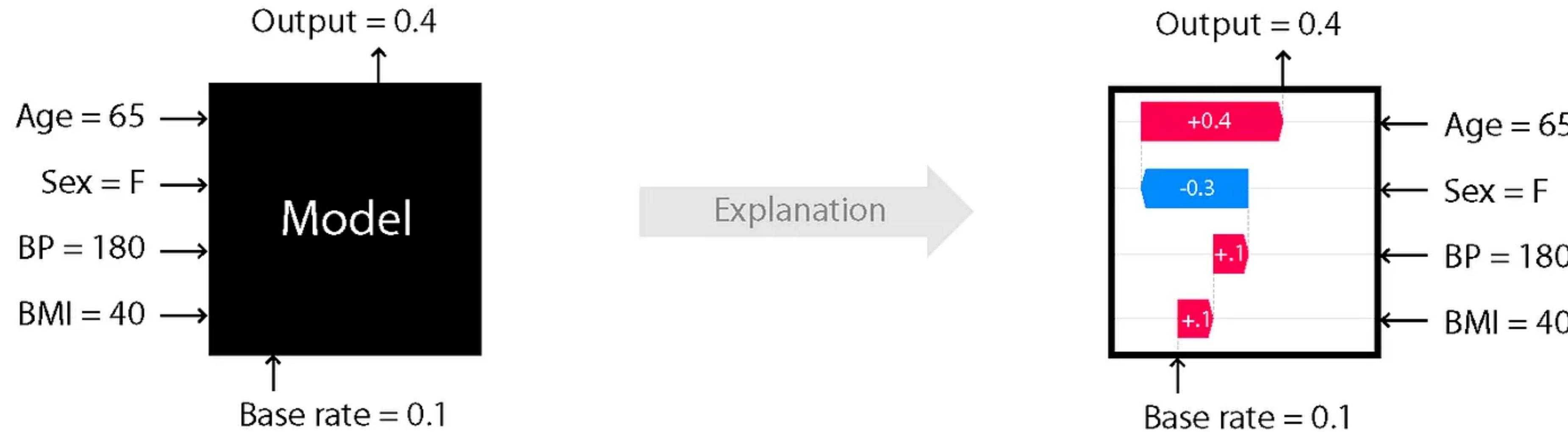
# we want to explain biological models over many sequences



# we want to explain biological models over many sequences

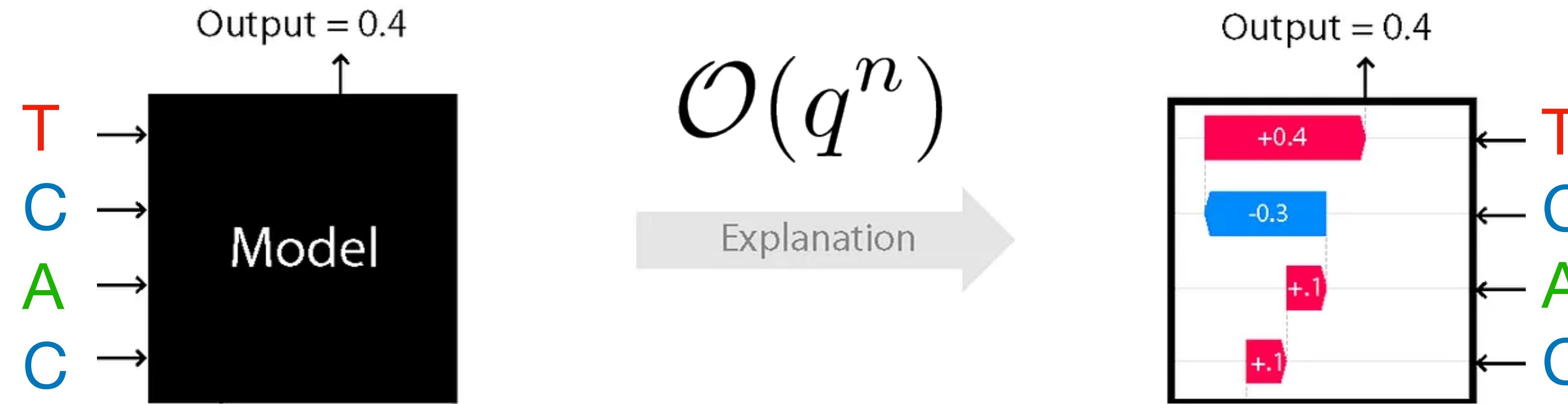


# we want to explain biological models over many sequences



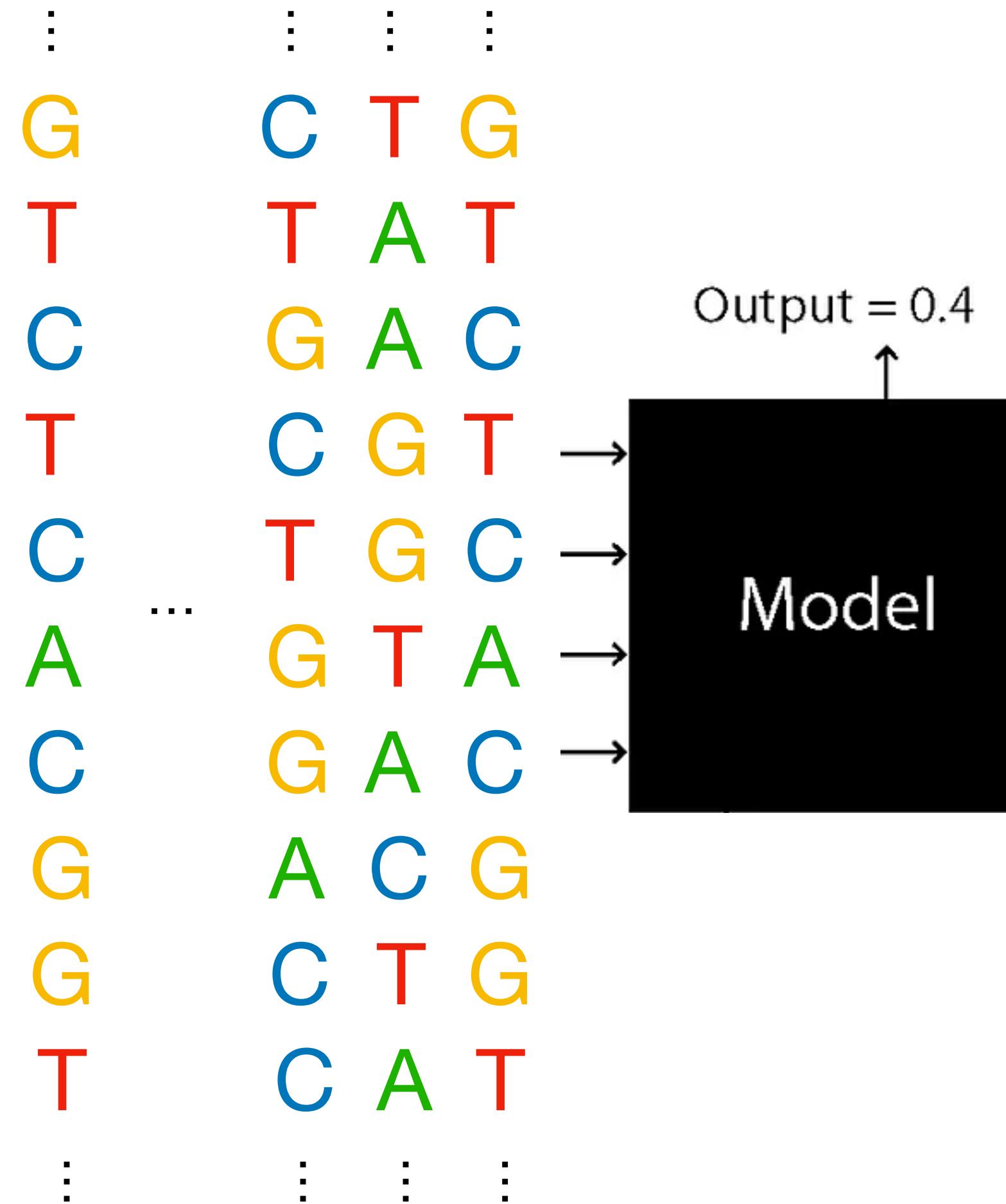
$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

# we want to explain biological models over many sequences



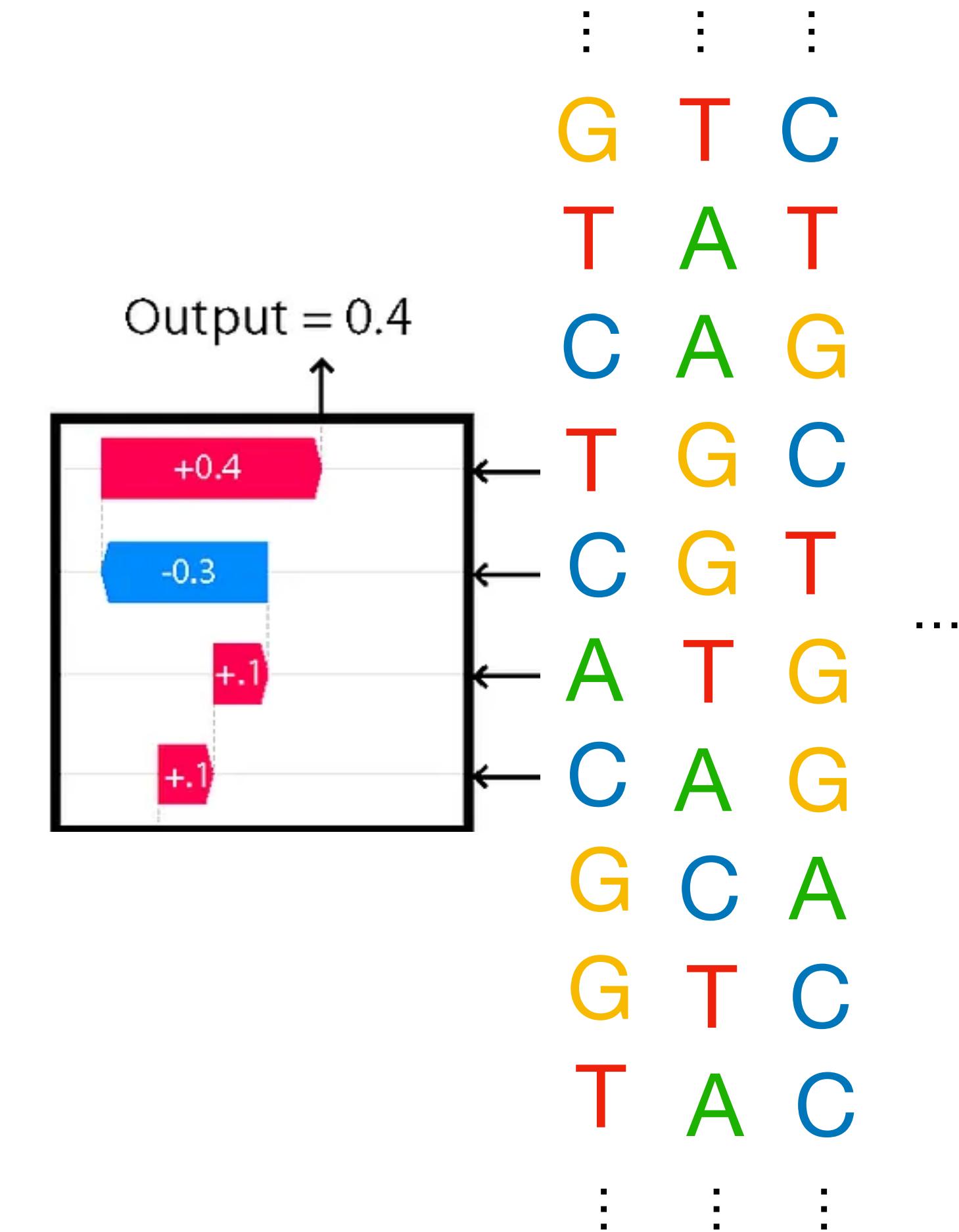
$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

# we want to explain biological models over many sequences



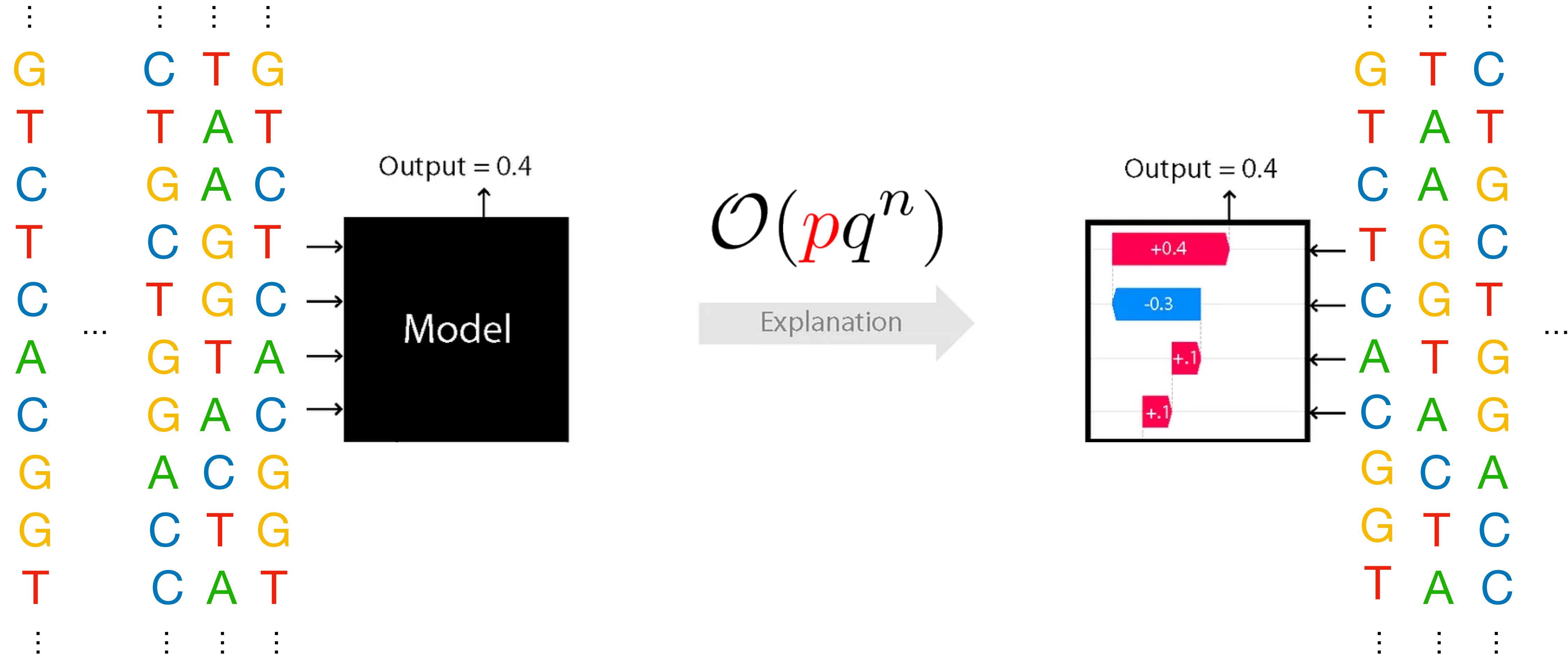
$\mathcal{O}(pq^n)$

Explanation →



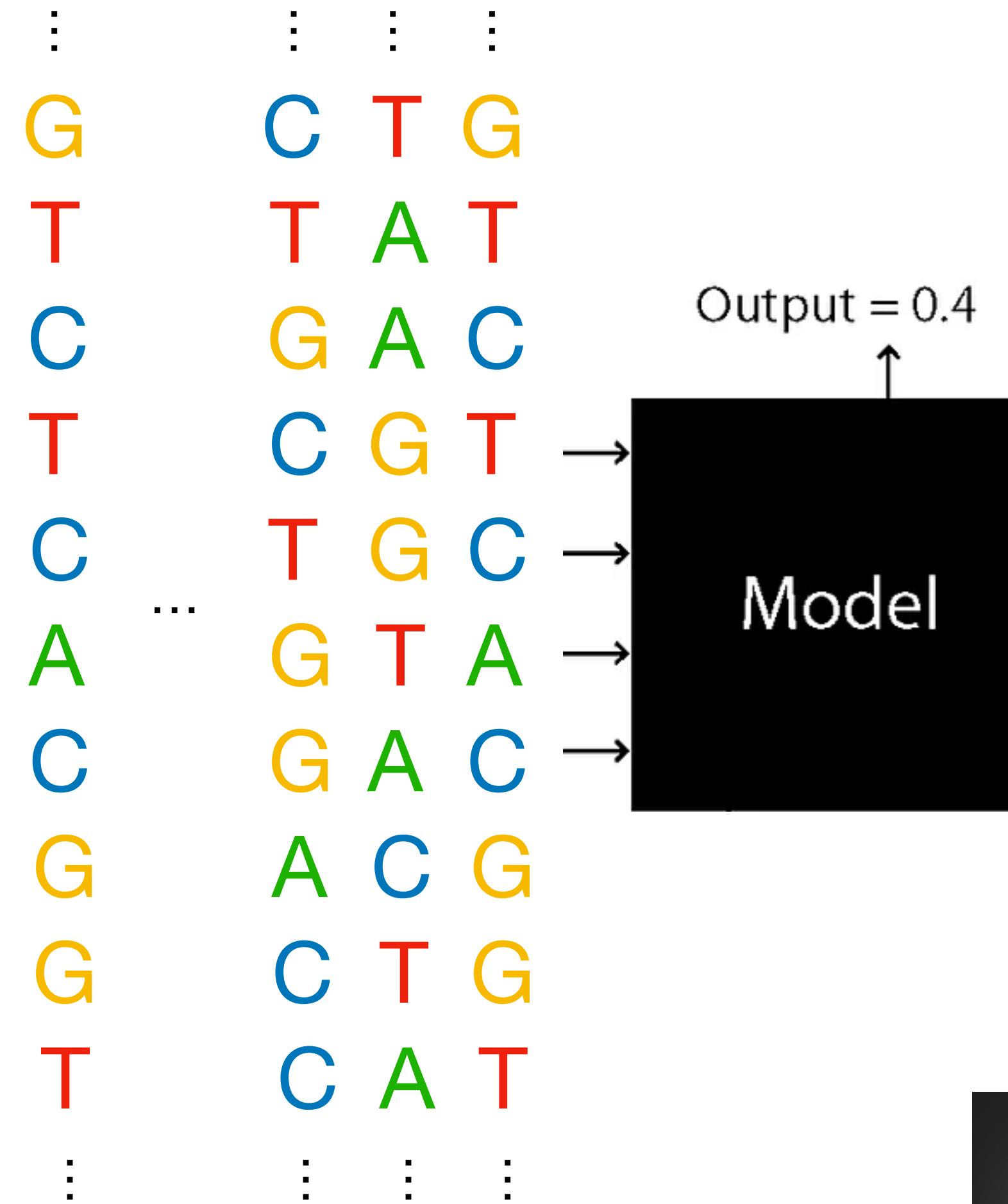
$$I_{\mathbf{x}_i}^{SV}(j) = \sum_{T \subseteq D \setminus \{j\}} \frac{|T|! (|D| - |T| - 1)!}{|D|!} [v_{T \cup \{j\}}(\mathbf{x}_i) - v_T(\mathbf{x}_i)]$$

# we want to explain biological models over many sequences



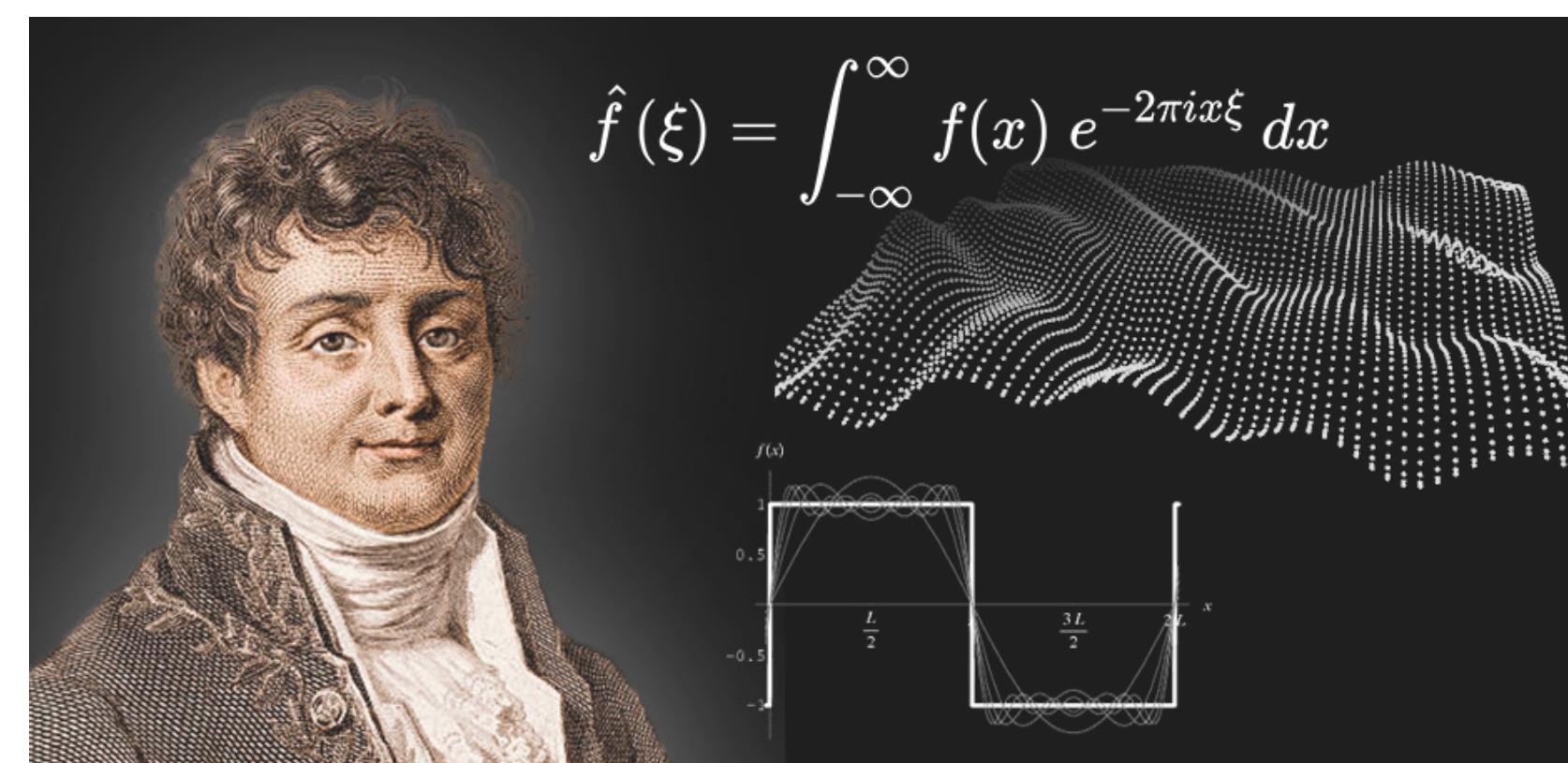
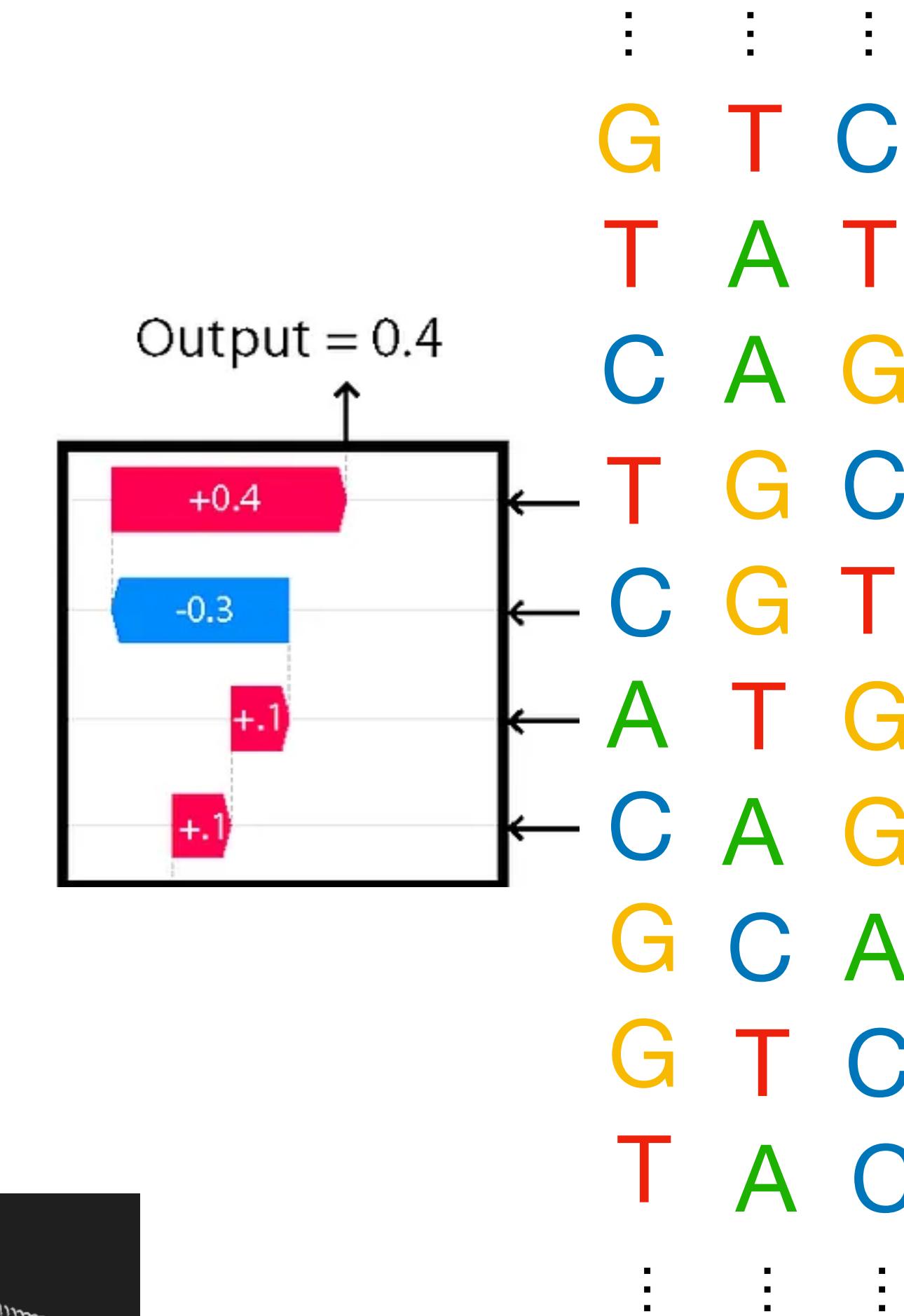
Explaining 1038 sequences  
Faith-Shap ~ 81 days!

# we want to explain biological models over many sequences



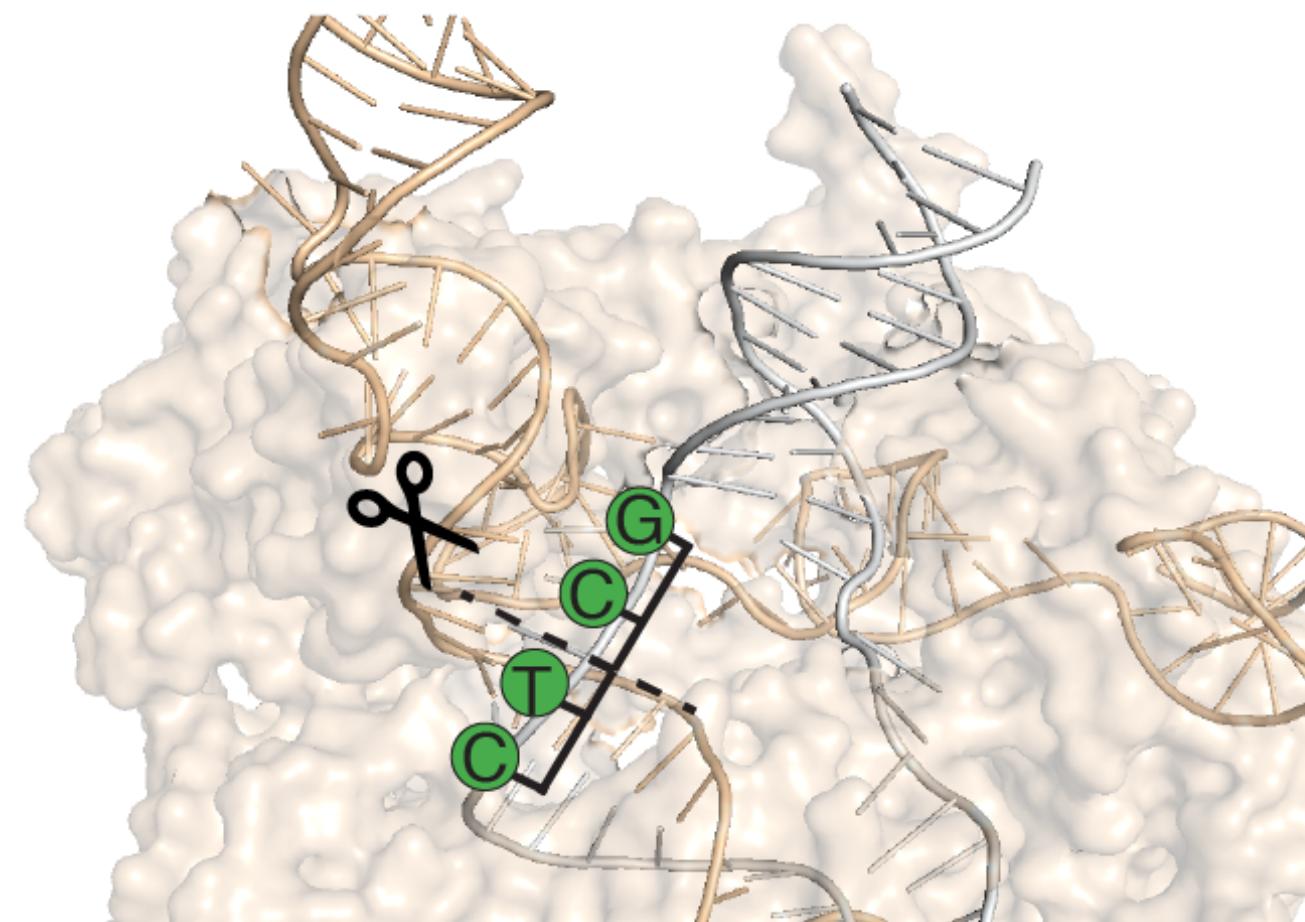
$$\mathcal{O}(pq^n)$$

Explanation



# biological sequence models are fundamentally sparse

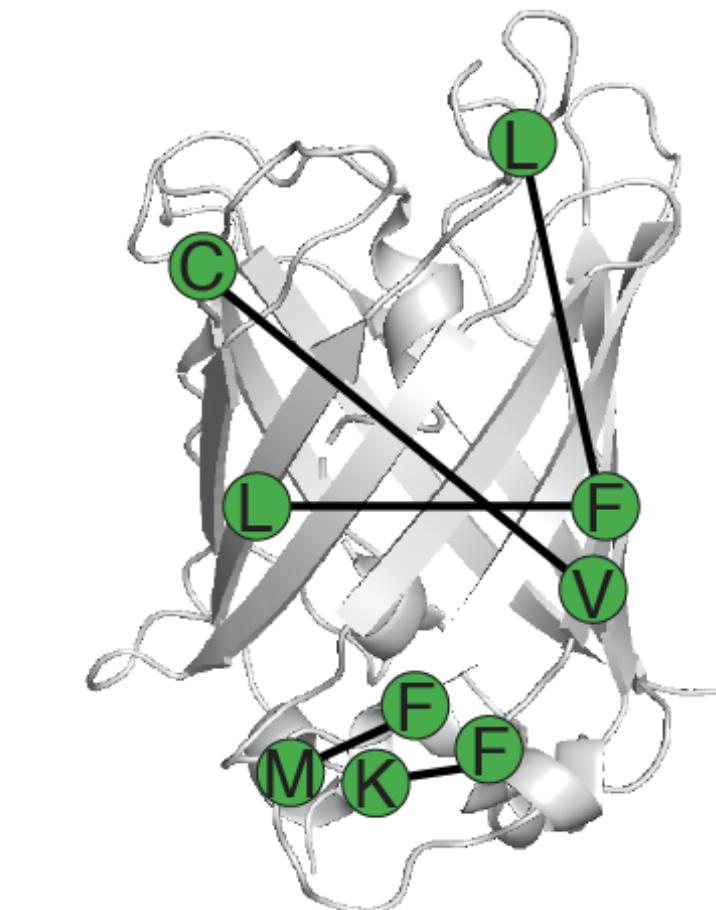
**a** Interactions in DNA repair



**b** Interactions in guide RNA binding

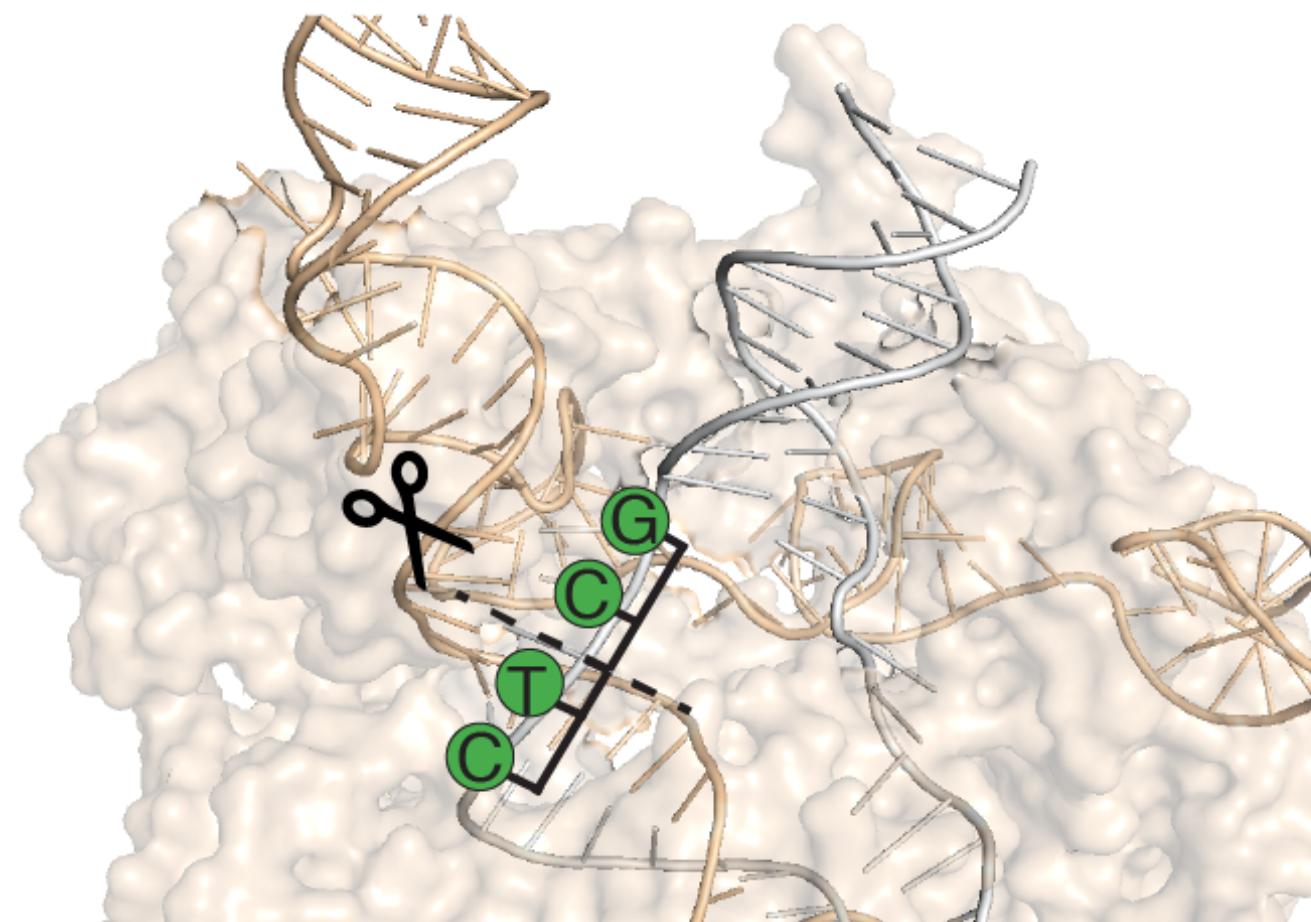


**c** Interactions in protein function



# biological sequence models are fundamentally sparse

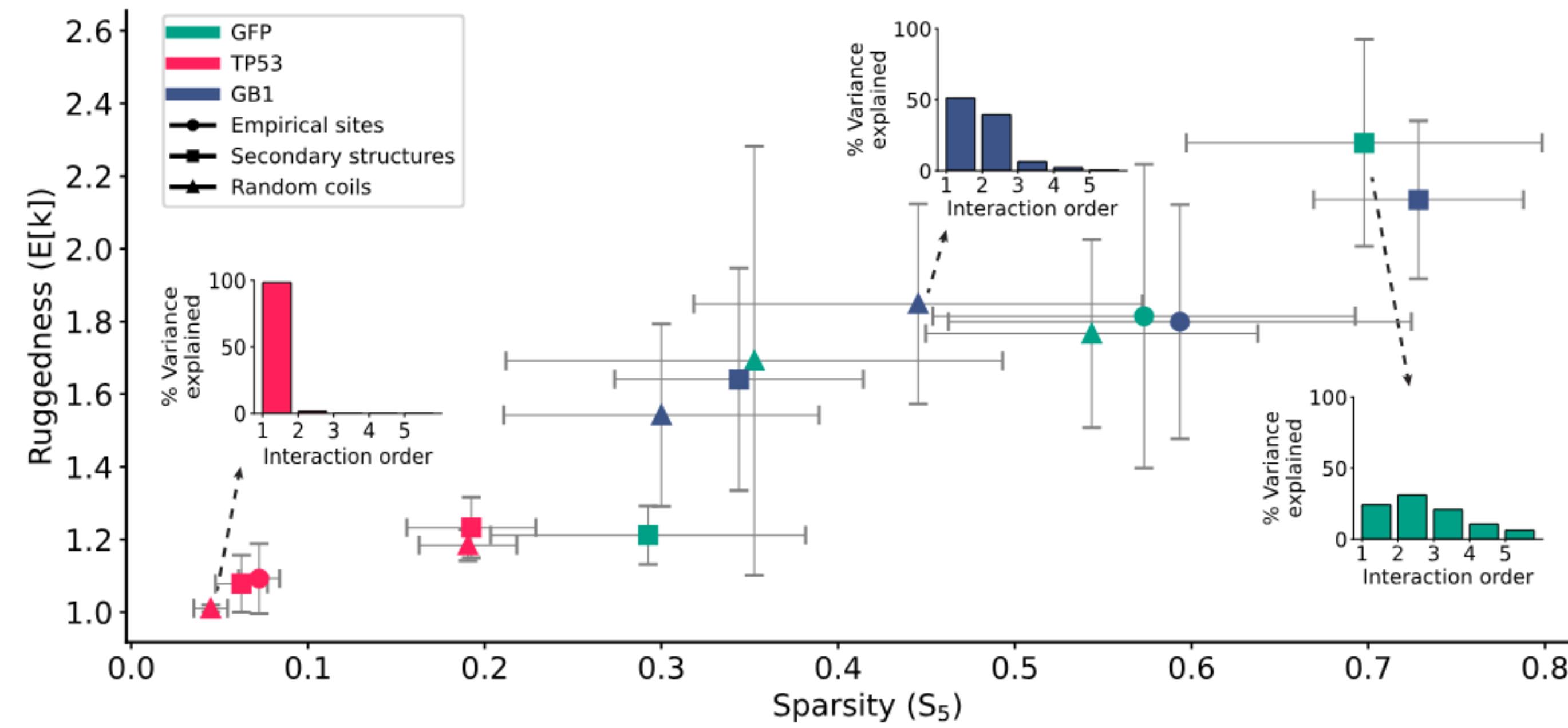
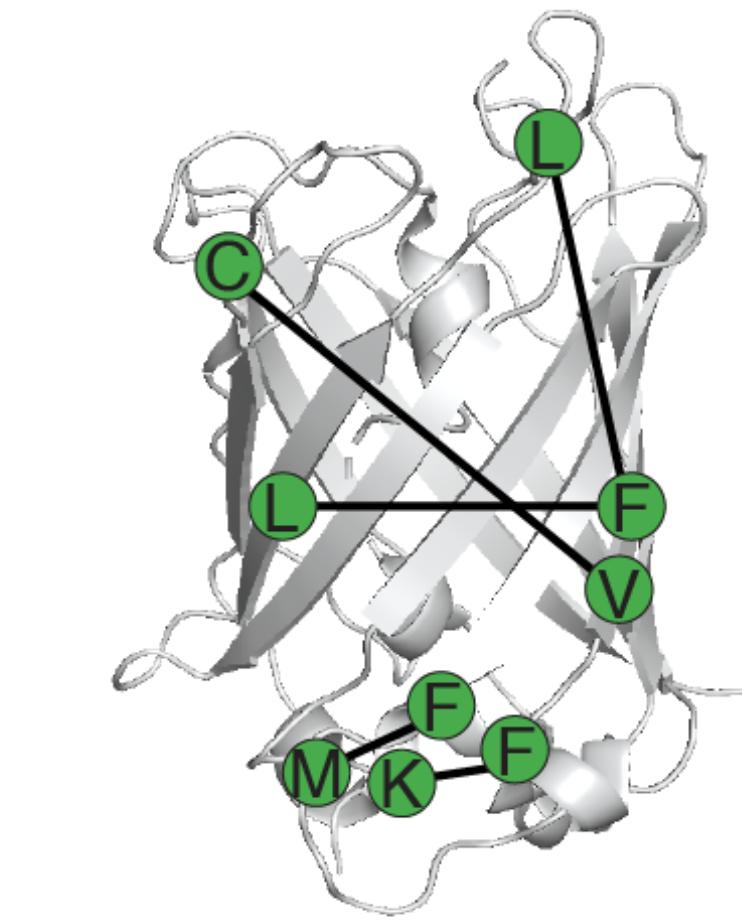
**a** Interactions in DNA repair



**b** Interactions in guide RNA binding



**c** Interactions in protein function



# biological sequence models as $q$ -ary functions

$$f : \mathbb{Z}_q^n$$

$$f([x_1, x_2, \dots, x_n]) \rightarrow y$$

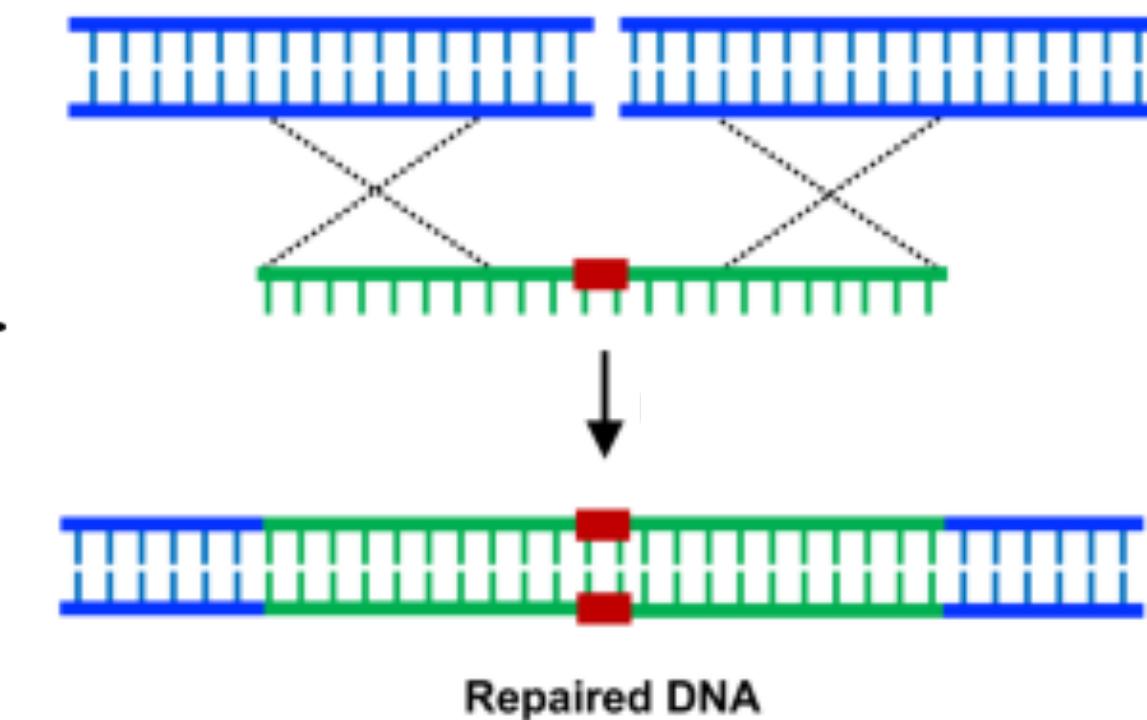
$$x_i \in \{0, 1, \dots, q - 1\}$$

# biological sequence models as $q$ -ary functions

$$f : \mathbb{Z}_4^n$$

$$f([x_1, x_2, \dots, x_n]) \rightarrow$$

$$x_i \in \{A, C, G, T\}$$



# biological sequence models as $q$ -ary functions

$$f : \mathbb{Z}_{20}^n$$

$$f([x_1, x_2, \dots, x_n]) \rightarrow$$



$$x_i \in \{P, L, \dots, Y\}$$

# biological sequence models as $q$ -ary functions

$$f : \mathbb{Z}_{20}^n$$

$$f([x_1, x_2, \dots, x_n]) \rightarrow$$



$$x_i \in \{P, L, \dots, Y\}$$

$$f(\mathbf{m}) = \sum_{\mathbf{y} \in \mathbb{Z}_q^n} F[\mathbf{y}] \omega^{\langle \mathbf{m}, \mathbf{y} \rangle} \quad \omega = e^{\frac{2\pi j}{q}}$$

# SHAP zero amortizes Shapley explanations

**Algorithm.** SHAP zero in three steps:

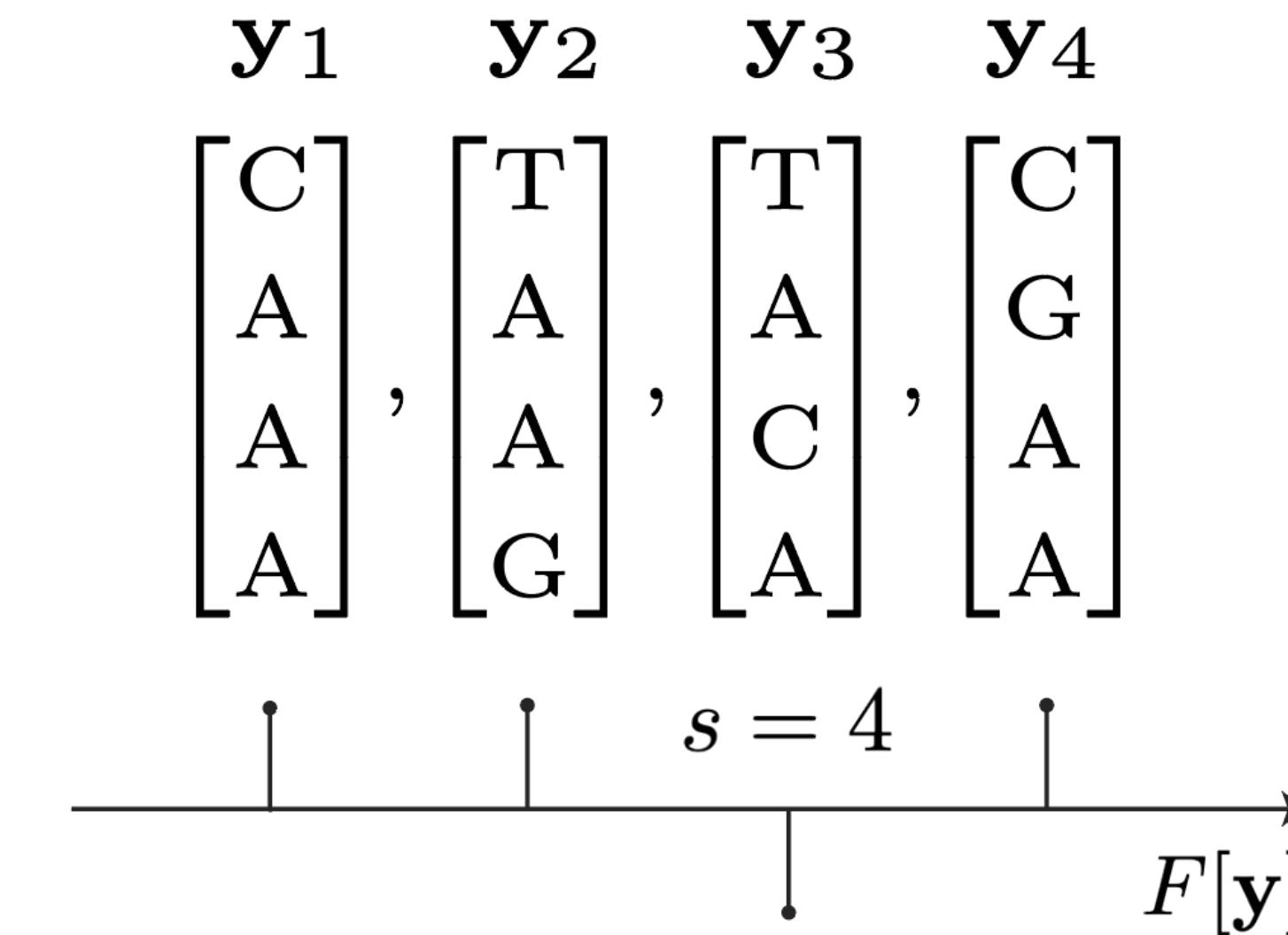
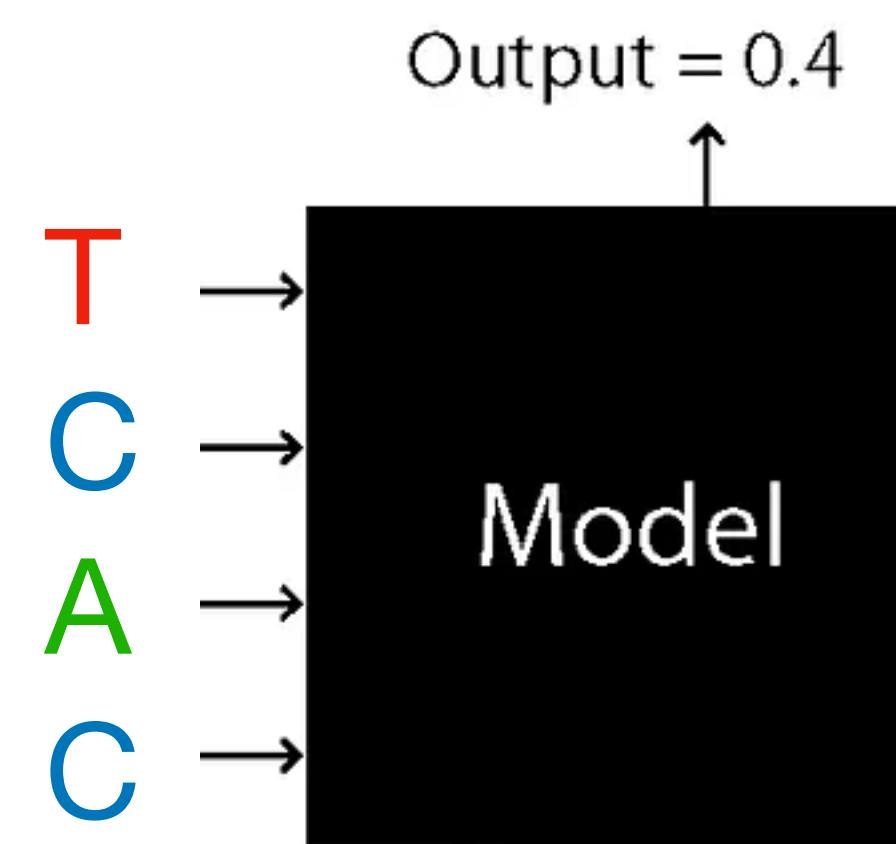
- **Step 1:** Pay a **one-time cost** to sketch the model globally via the Fourier transform with **sample complexity**  $\mathcal{O}(sn^2)$  and **computational complexity**  $\mathcal{O}(sn^3)$

**For each future sequence:**

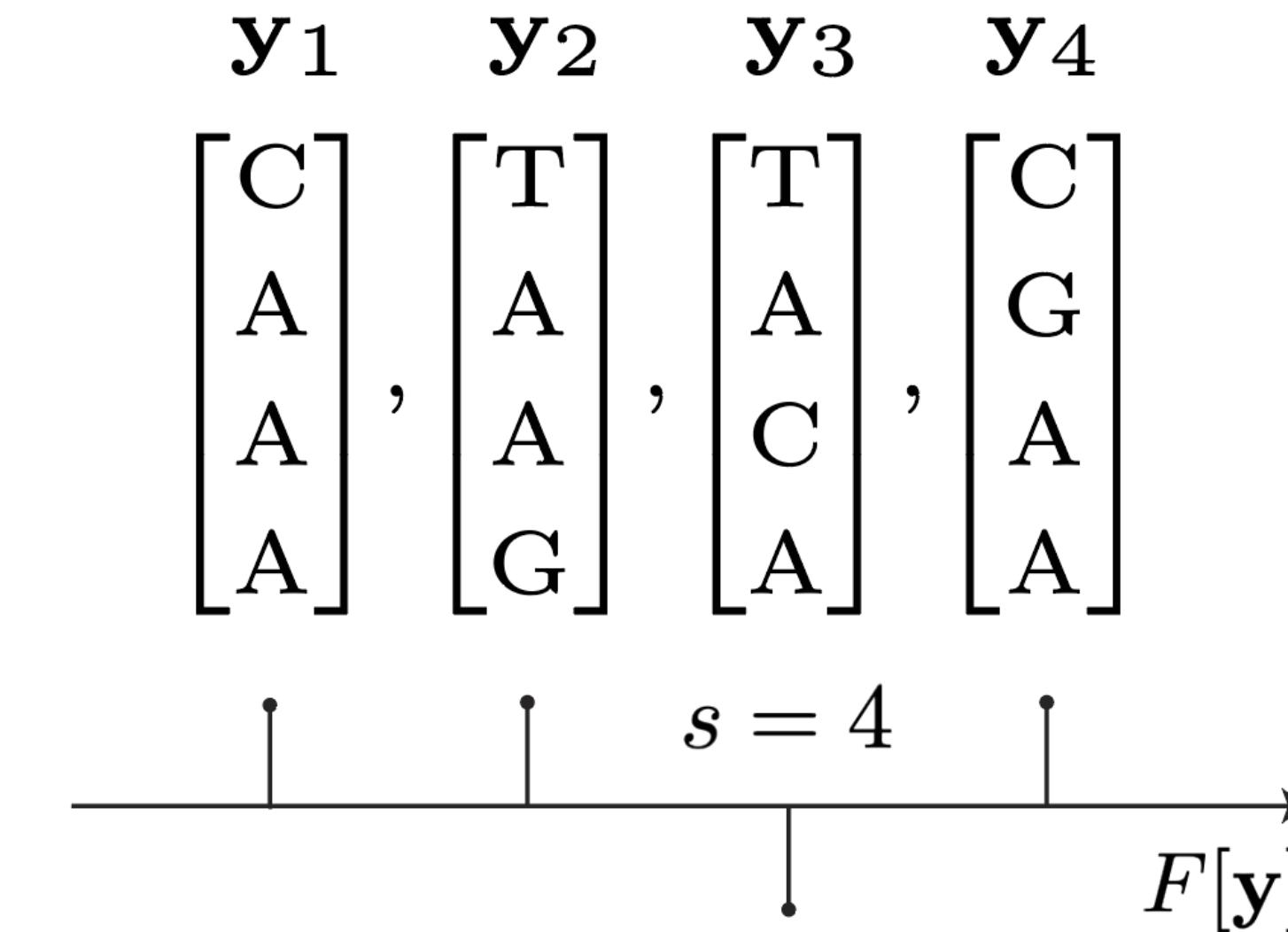
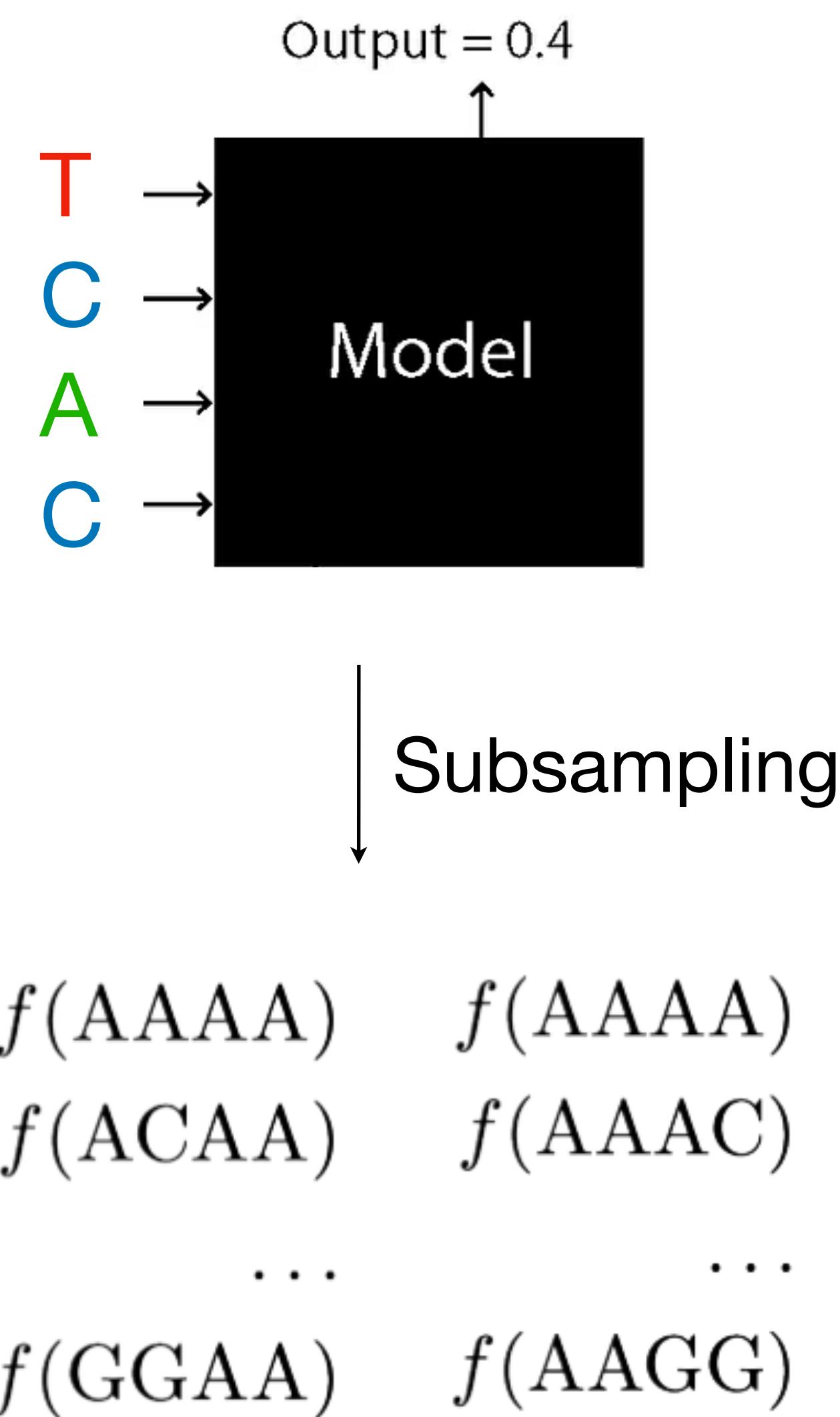
- **Step 2:** Map global sketch to the local Möbius transform to isolate interacting features with complexity  $\mathcal{O}(s^2(2q)^\ell)$
- **Step 3:** Map Möbius transform to Shapley values and interactions with complexity  $\mathcal{O}(s^2(2q)^\ell)$



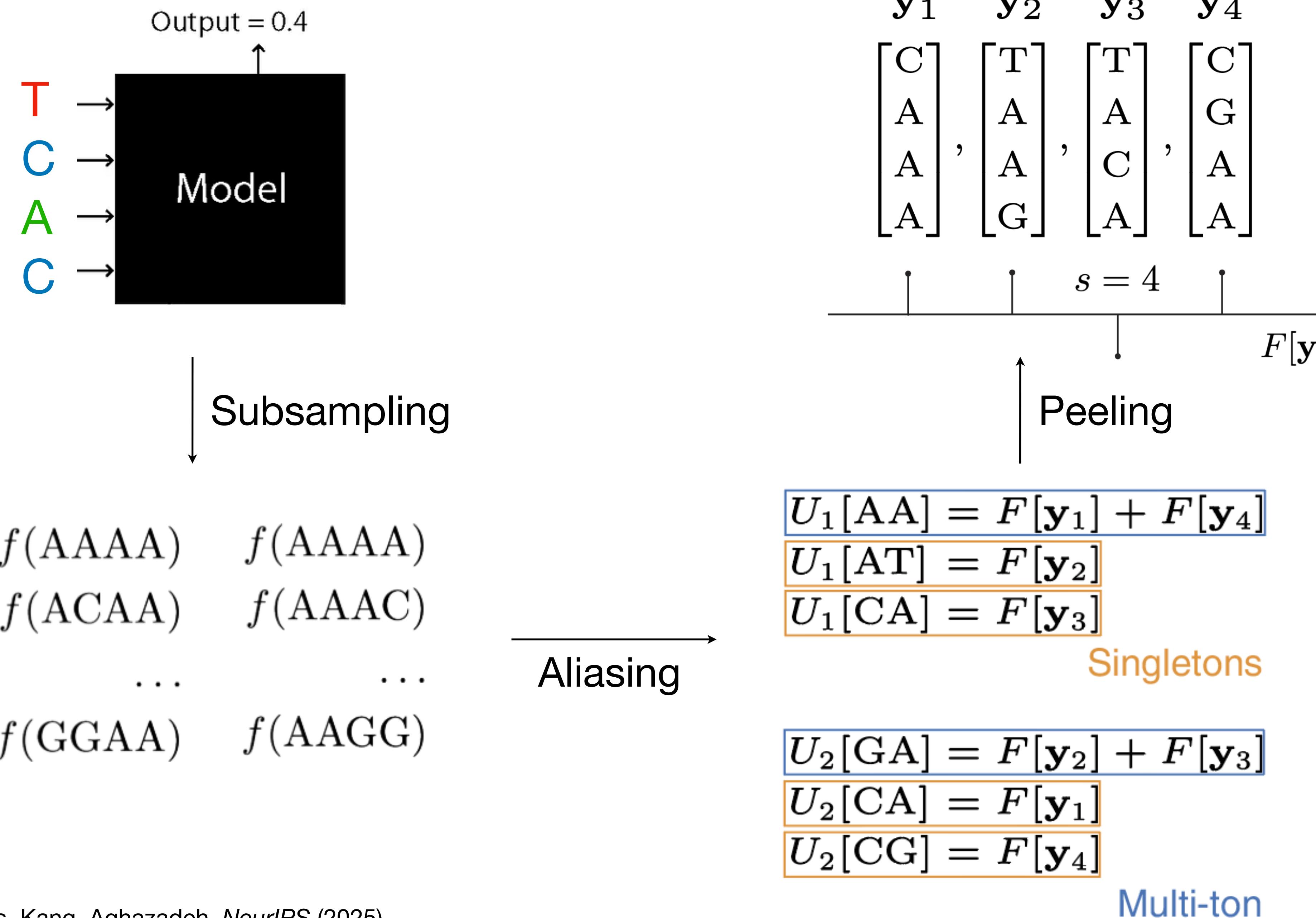
# step 1: pay a **one-time cost** to sketch the model globally



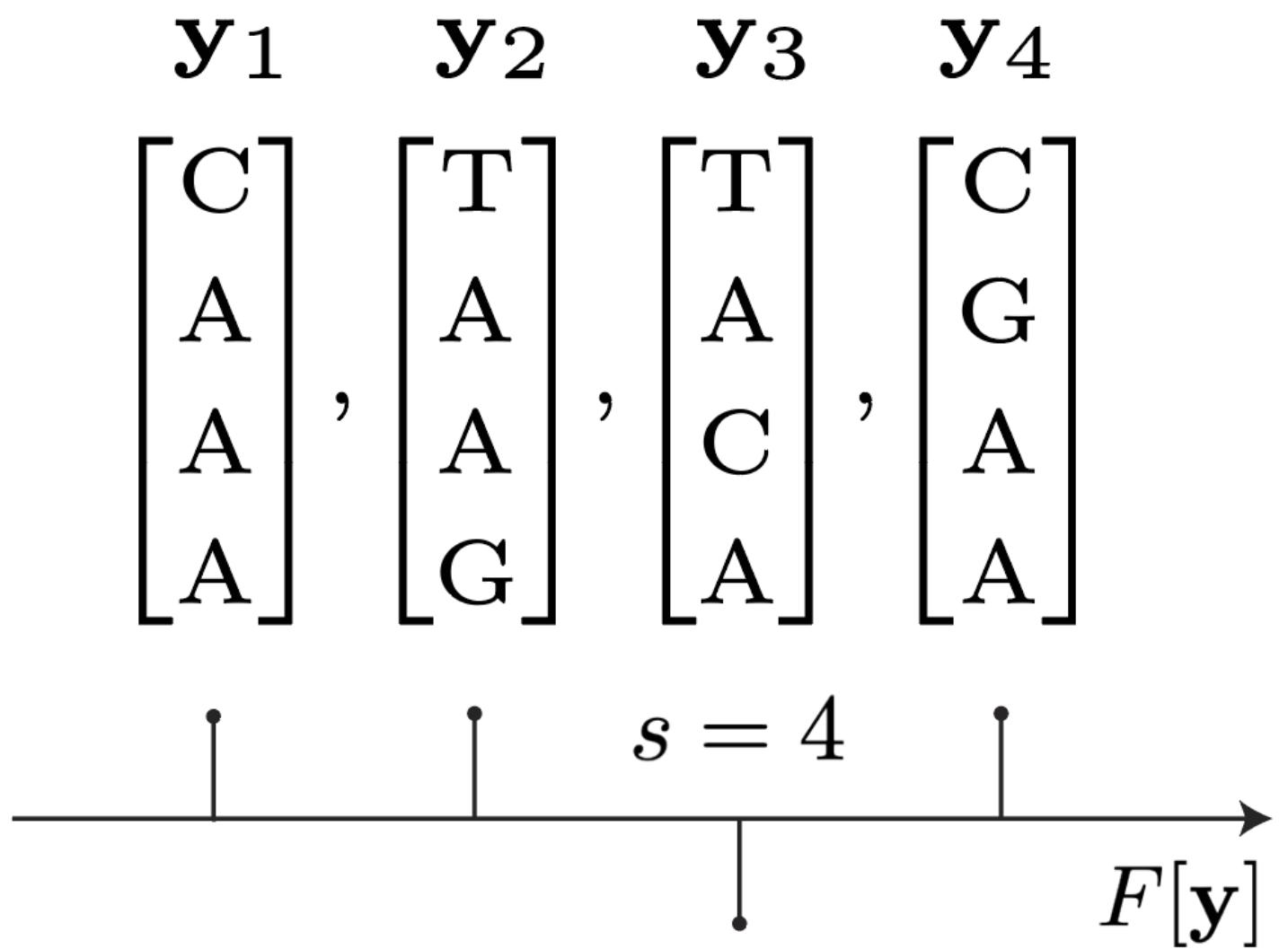
# step 1: pay a **one-time cost** to sketch the model globally



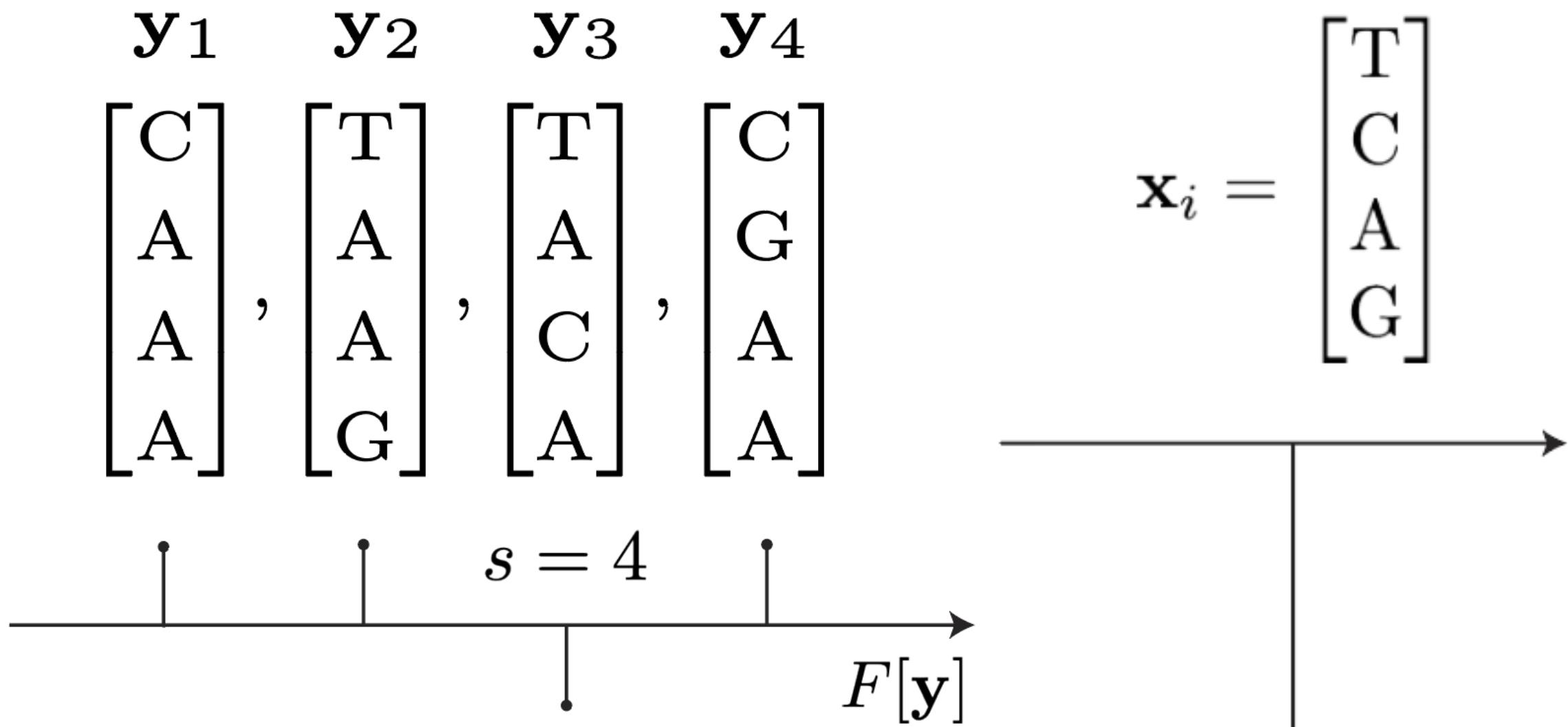
# step 1: pay a one-time cost to sketch the model globally



## step 2: map sketch to local Möbius transform



## step 2: map sketch to local Möbius transform

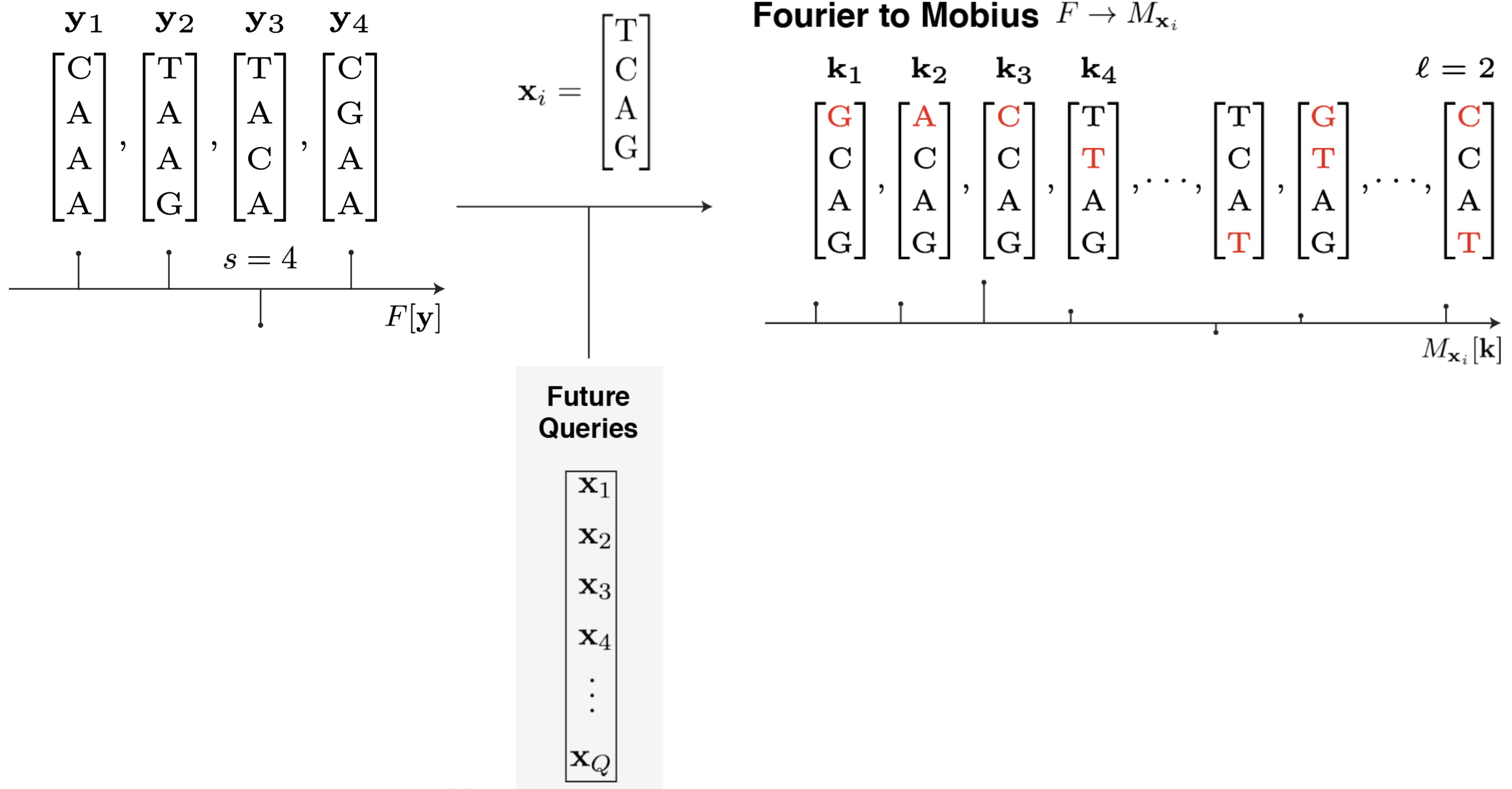


$$\mathbf{x}_i = \begin{bmatrix} T \\ C \\ A \\ G \end{bmatrix}$$

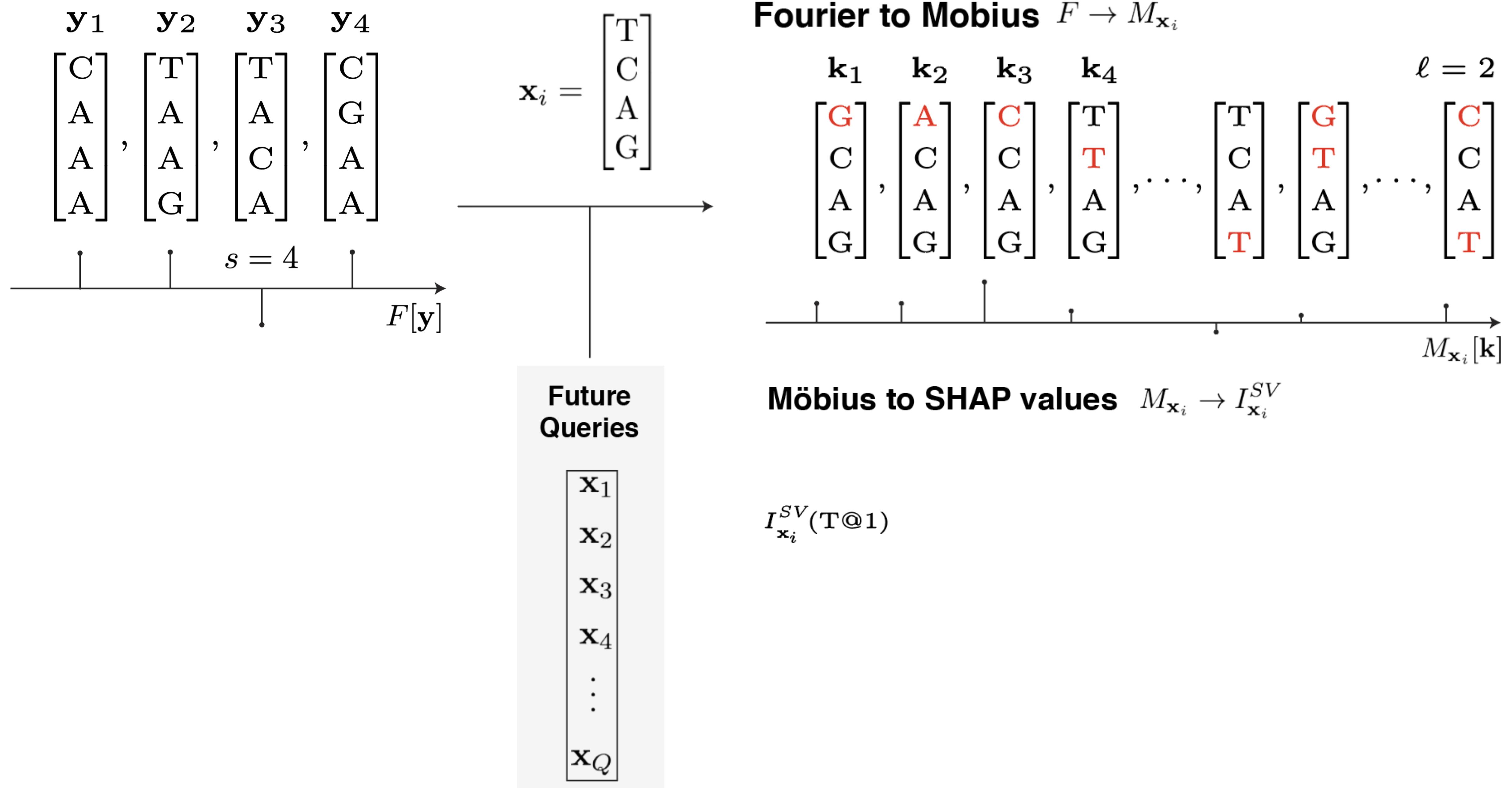
Future  
Queries

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \vdots \\ \mathbf{x}_Q \end{bmatrix}$$

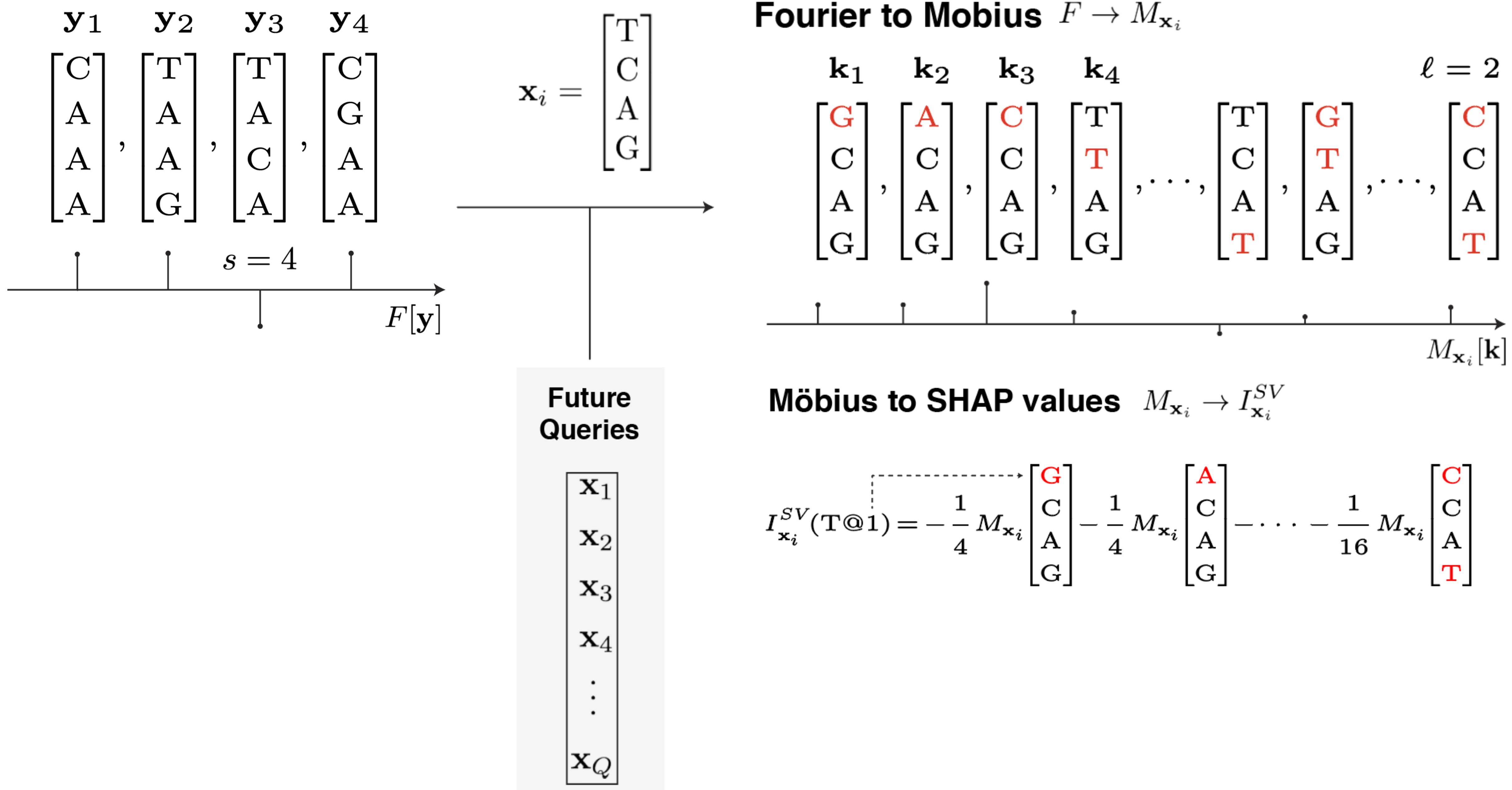
## step 2: map sketch to local Möbius transform



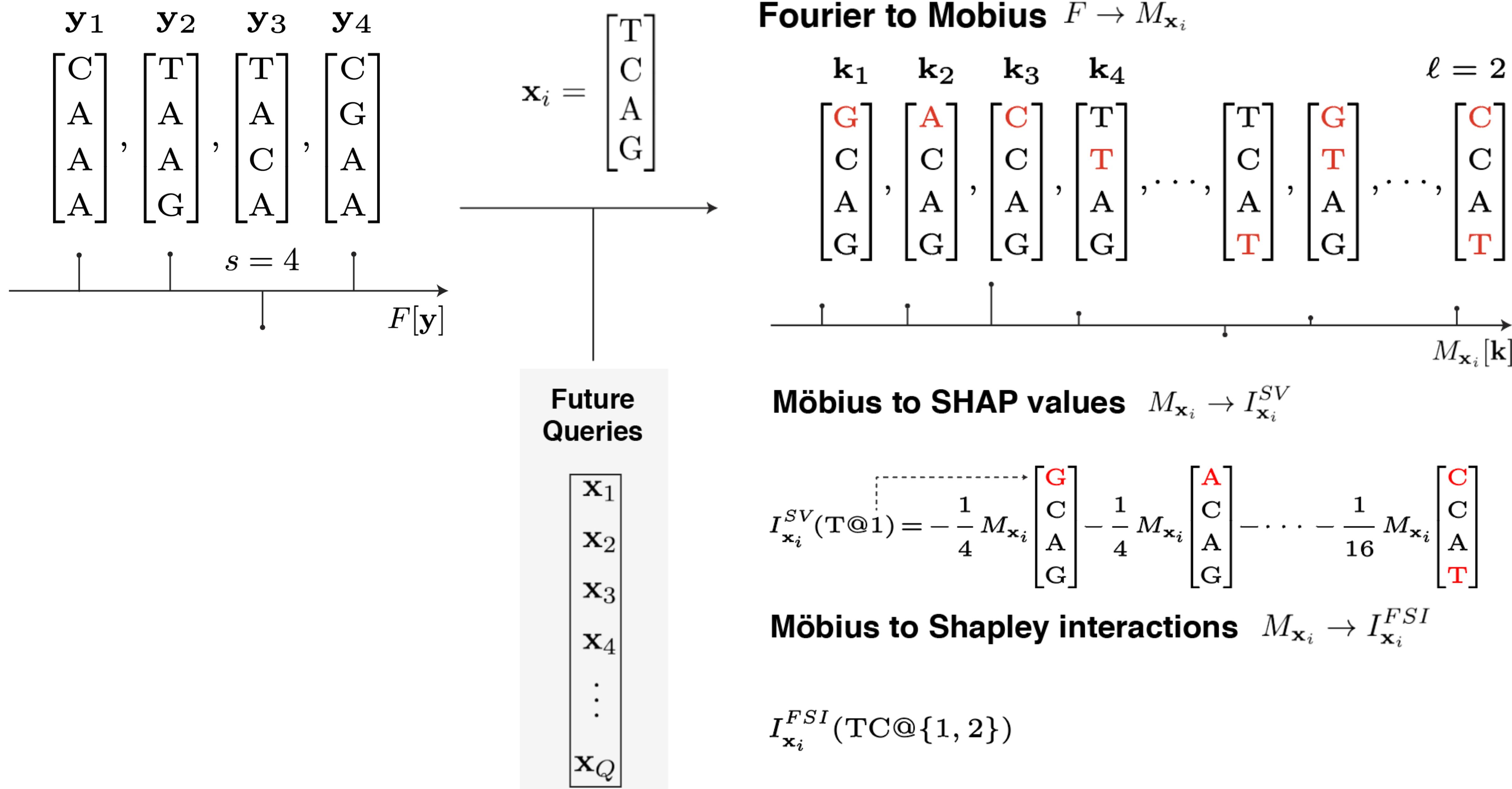
# step 3: map Möbius transform to SHAP values and interactions



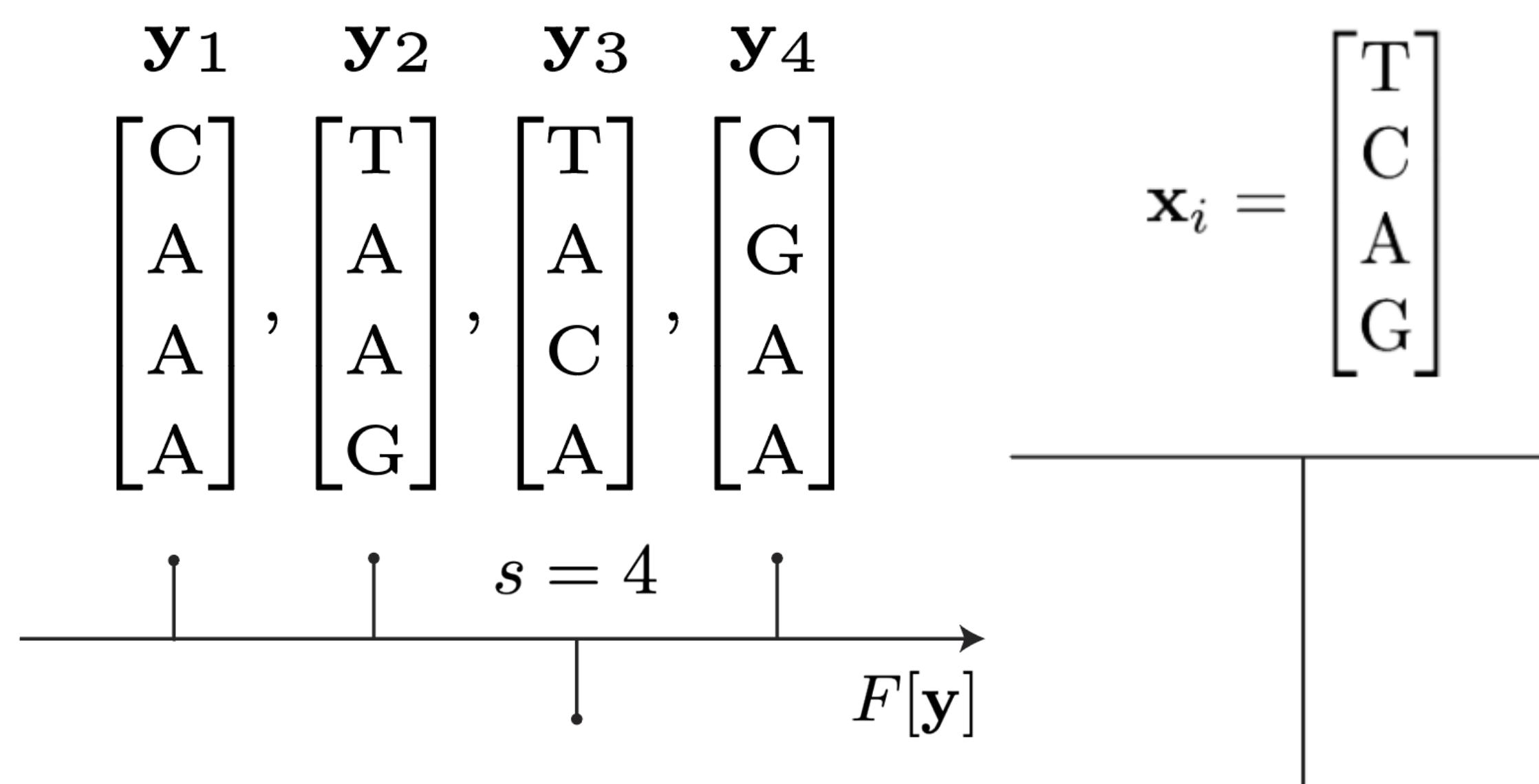
# step 3: map Möbius transform to SHAP values and interactions



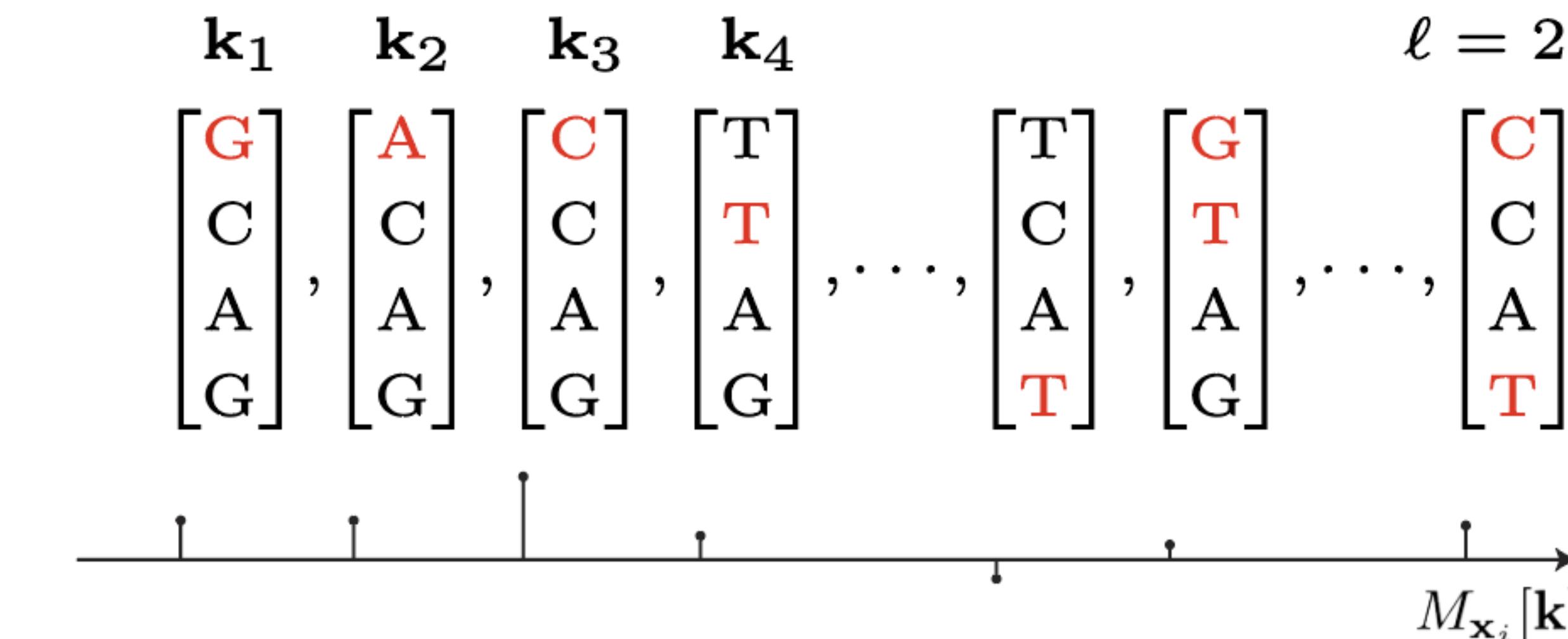
# step 3: map Möbius transform to SHAP values and interactions



# step 3: map Möbius transform to SHAP values and interactions



**Fourier to Möbius**  $F \rightarrow M_{\mathbf{x}_i}$



Future  
Queries

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \vdots \\ \mathbf{x}_Q \end{bmatrix}$$

**Möbius to SHAP values**  $M_{\mathbf{x}_i} \rightarrow I_{\mathbf{x}_i}^{SV}$

$$I_{\mathbf{x}_i}^{SV}(T@1) = -\frac{1}{4} M_{\mathbf{x}_i} \begin{bmatrix} G \\ C \\ A \\ G \end{bmatrix} - \frac{1}{4} M_{\mathbf{x}_i} \begin{bmatrix} A \\ C \\ A \\ G \end{bmatrix} - \dots - \frac{1}{16} M_{\mathbf{x}_i} \begin{bmatrix} C \\ C \\ A \\ T \end{bmatrix}$$

**Möbius to Shapley interactions**  $M_{\mathbf{x}_i} \rightarrow I_{\mathbf{x}_i}^{FSI}$

$$I_{\mathbf{x}_i}^{FSI}(TC@\{1, 2\}) = \frac{1}{16} M_{\mathbf{x}_i} \begin{bmatrix} G \\ T \\ A \\ G \end{bmatrix} + \dots + \frac{1}{16} M_{\mathbf{x}_i} \begin{bmatrix} C \\ G \\ A \\ G \end{bmatrix}$$

# SHAP zero amortizes Shapley explanations

**Algorithm.** SHAP zero in three steps:

- **Step 1:** Pay a **one-time cost** to sketch the model globally via the Fourier transform with **sample complexity**  $\mathcal{O}(sn^2)$  and **computational complexity**  $\mathcal{O}(sn^3)$

**For each future sequence:**

- **Step 2:** Map global sketch to the local Möbius transform to isolate interacting features with complexity  $\mathcal{O}(s^2(2q)^\ell)$
- **Step 3:** Map Möbius transform to Shapley values and interactions with complexity  $\mathcal{O}(s^2(2q)^\ell)$



# SHAP zero amortizes Shapley explanations

**Algorithm.** SHAP zero in three steps:

- **Step 1:** Pay a **one-time cost** to sketch the model globally via the Fourier transform with **sample complexity**  $\mathcal{O}(sn^2)$  and **computational complexity**  $\mathcal{O}(sn^3)$

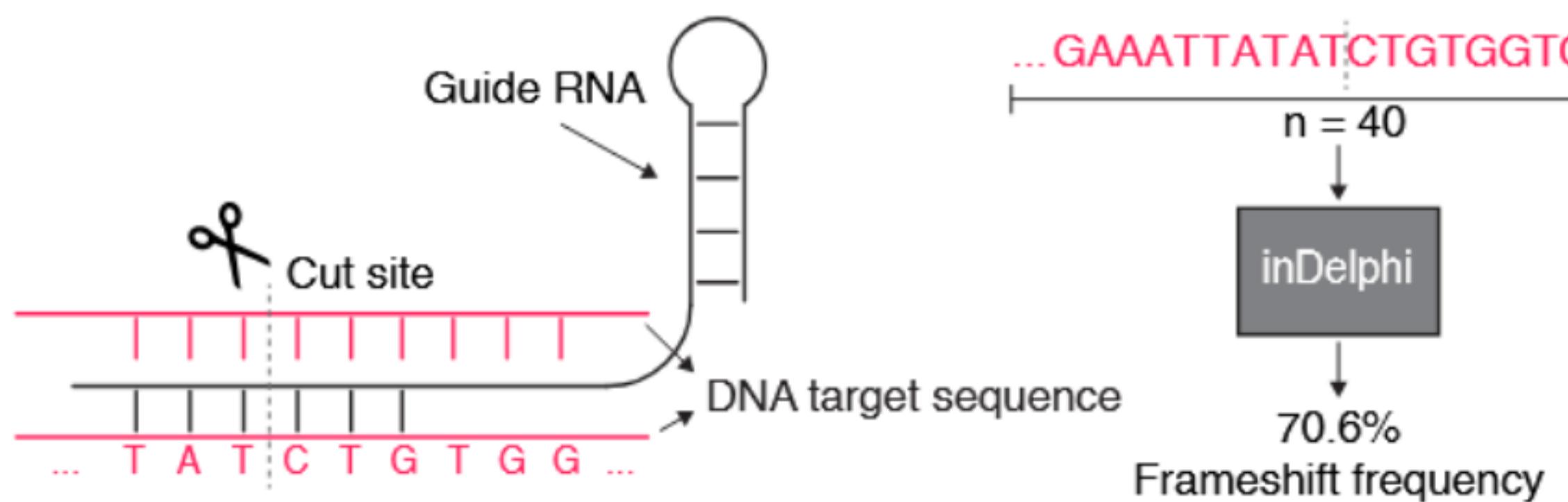
**For each future sequence:**

- **Step 2:** Map global sketch to the local Möbius transform to isolate interacting features with complexity  $\mathcal{O}(s^2(2q)^\ell)$
- **Step 3:** Map Möbius transform to Shapley values and interactions with complexity  $\mathcal{O}(s^2(2q)^\ell)$

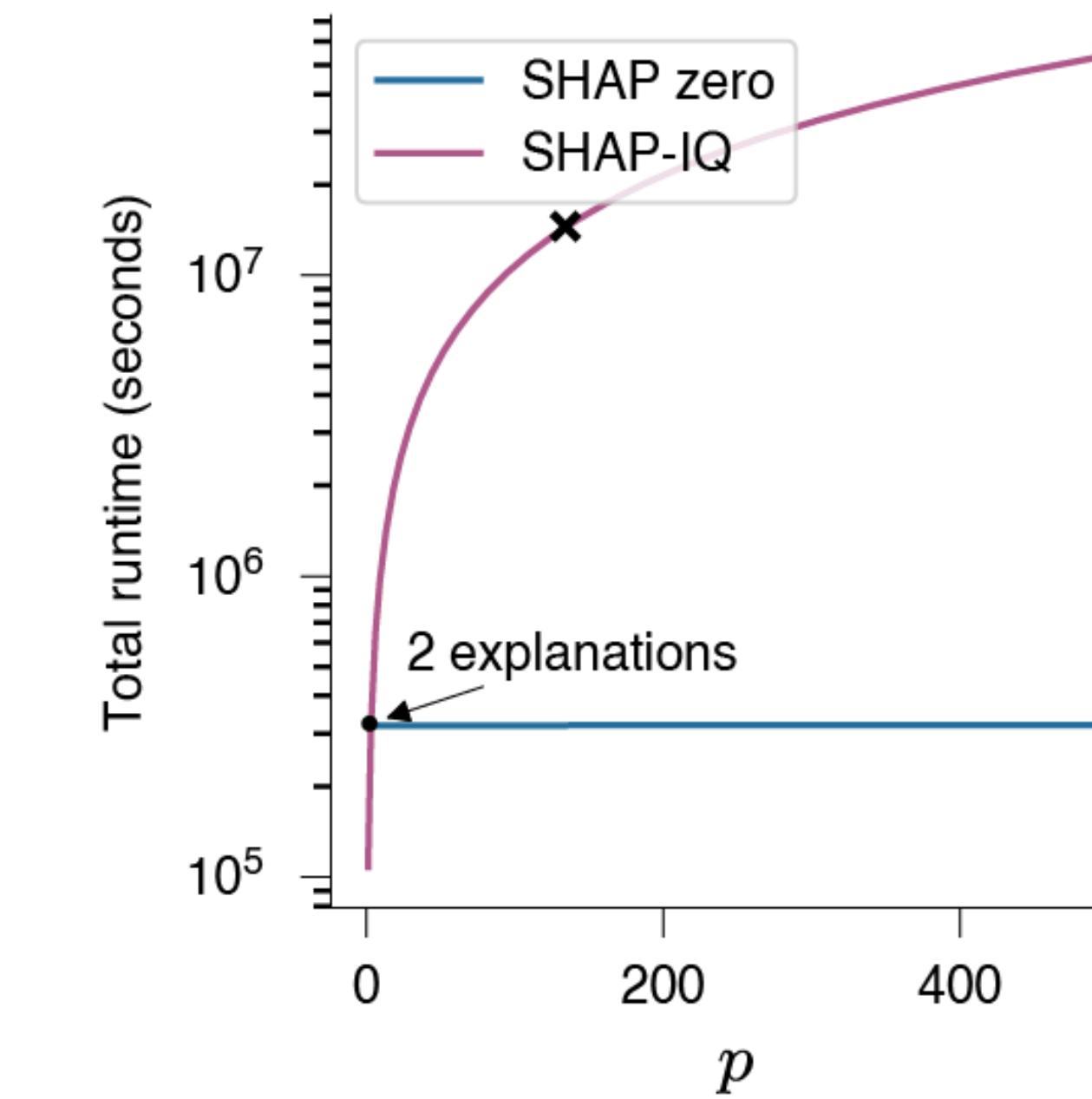
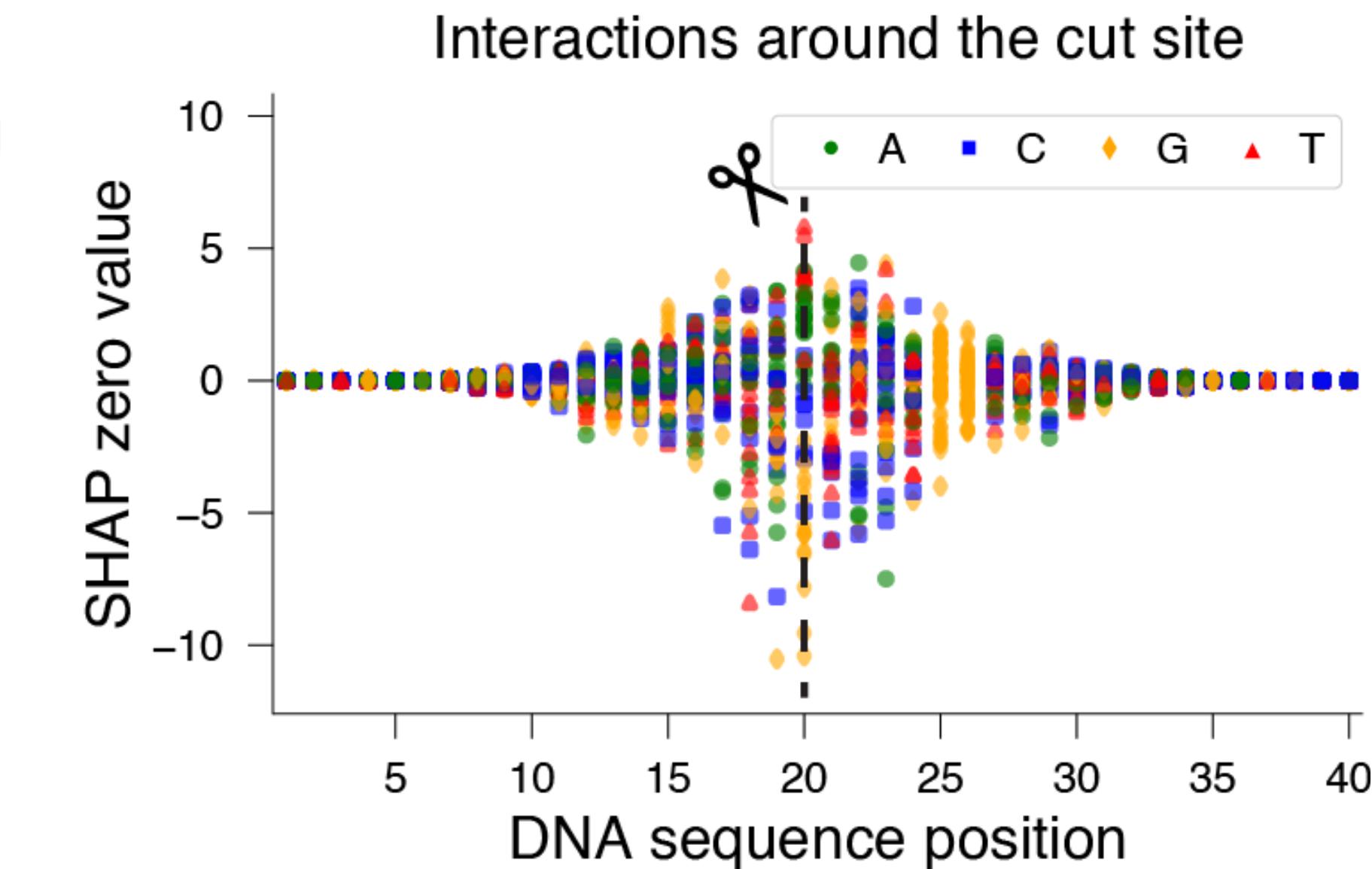
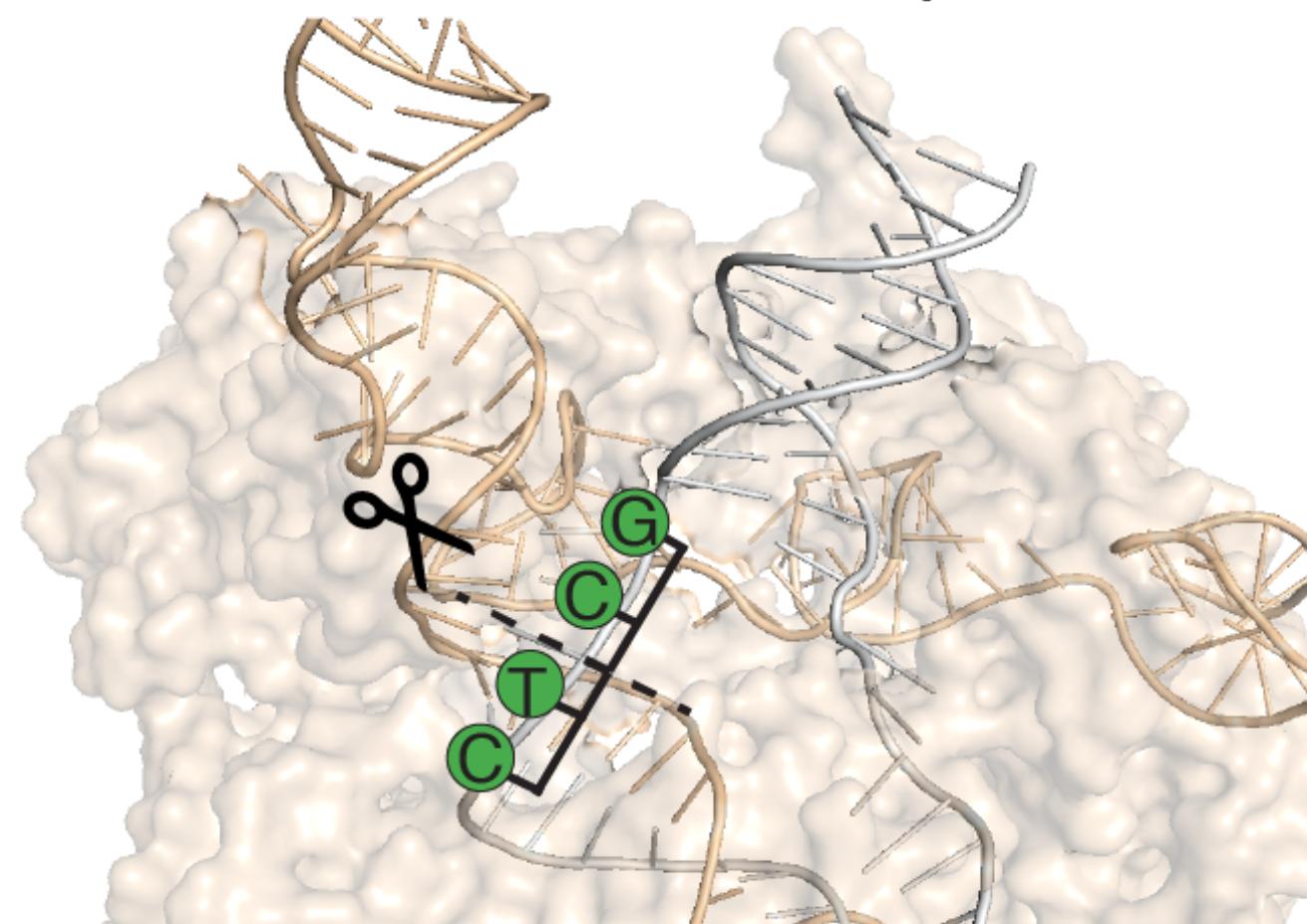


After Step 1, SHAP zero scales **essentially free!**

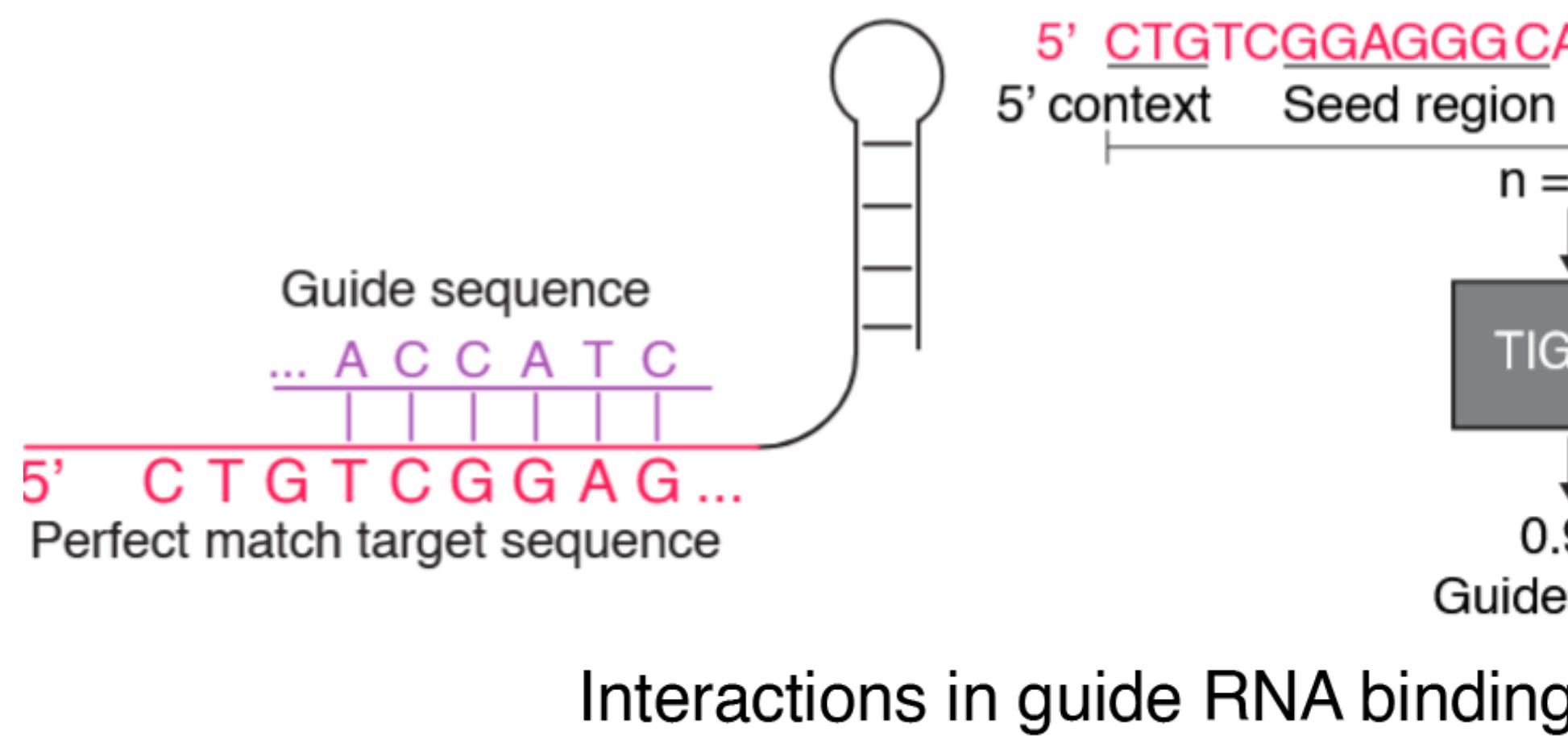
# SHAP zero uncovers DNA repair interactions at scale



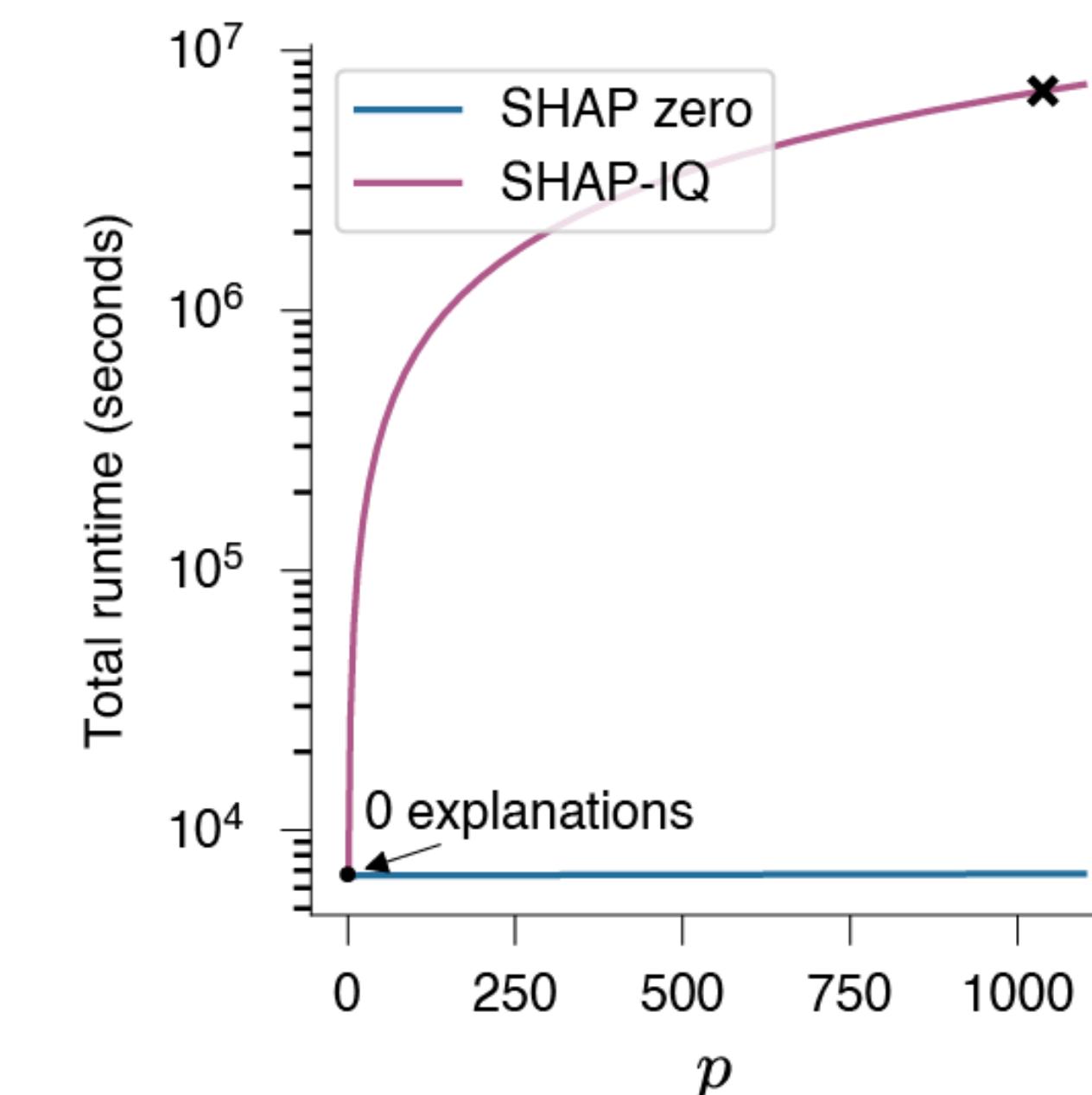
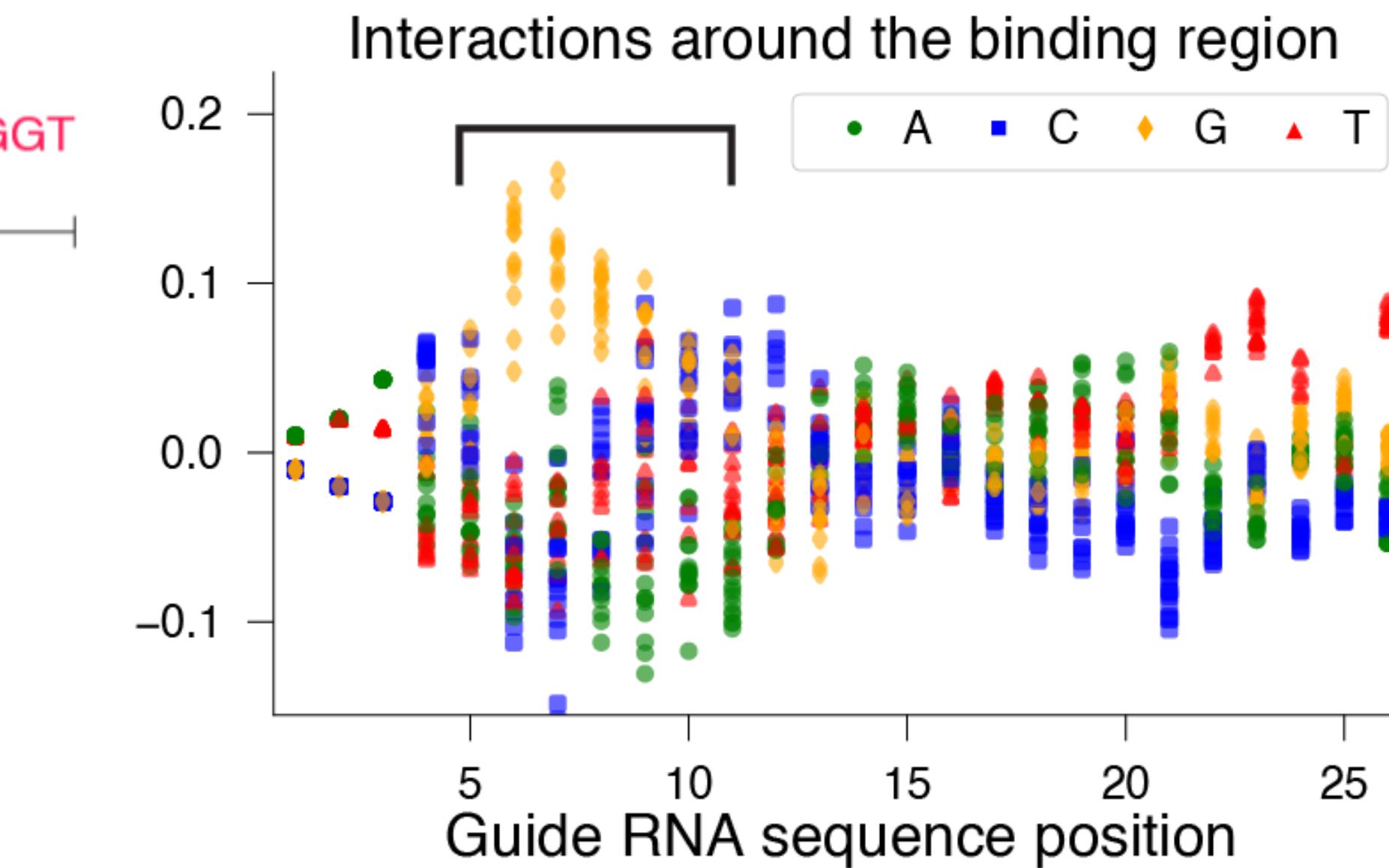
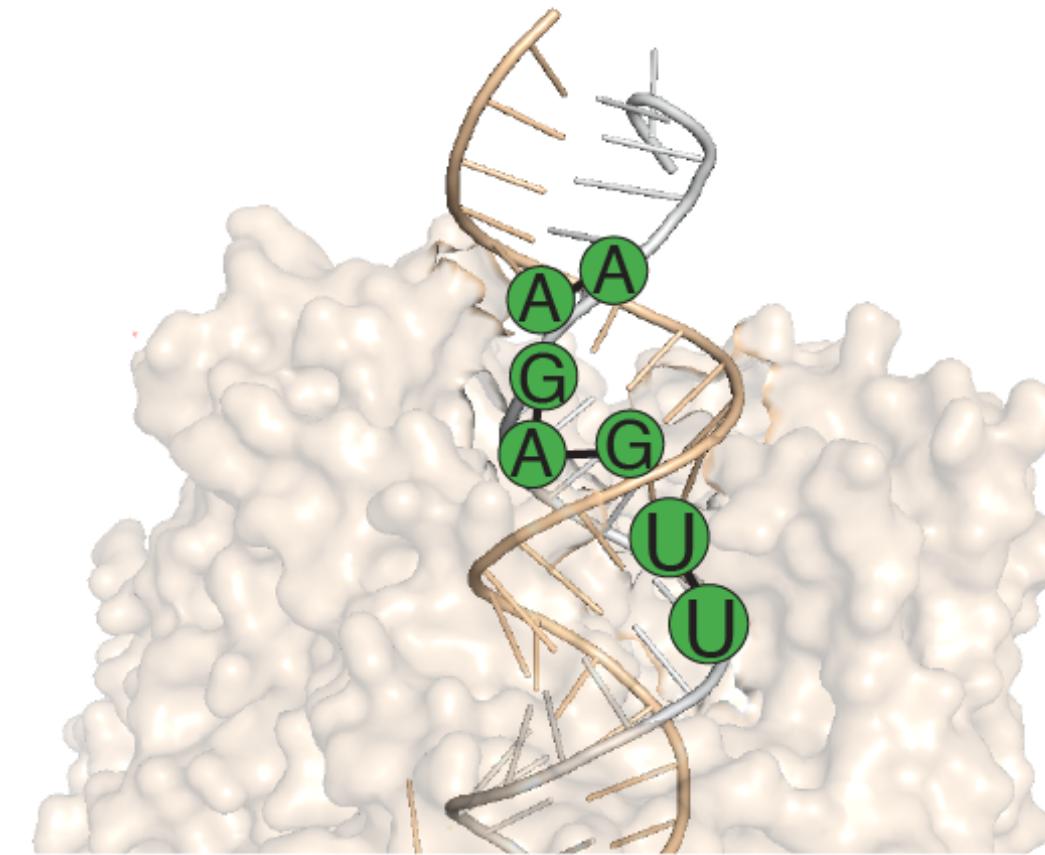
Interactions in DNA repair



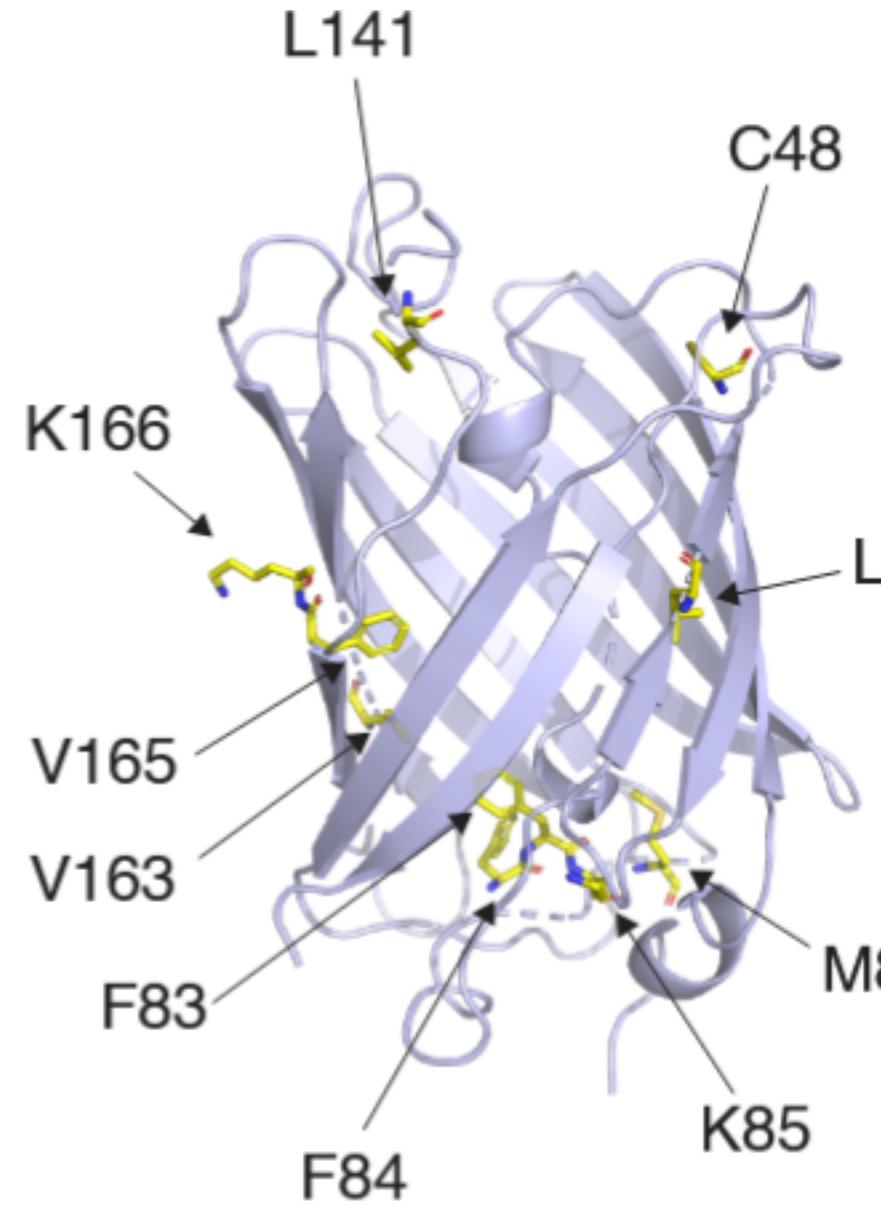
# SHAP zero uncovers guide RNA binding interactions at scale



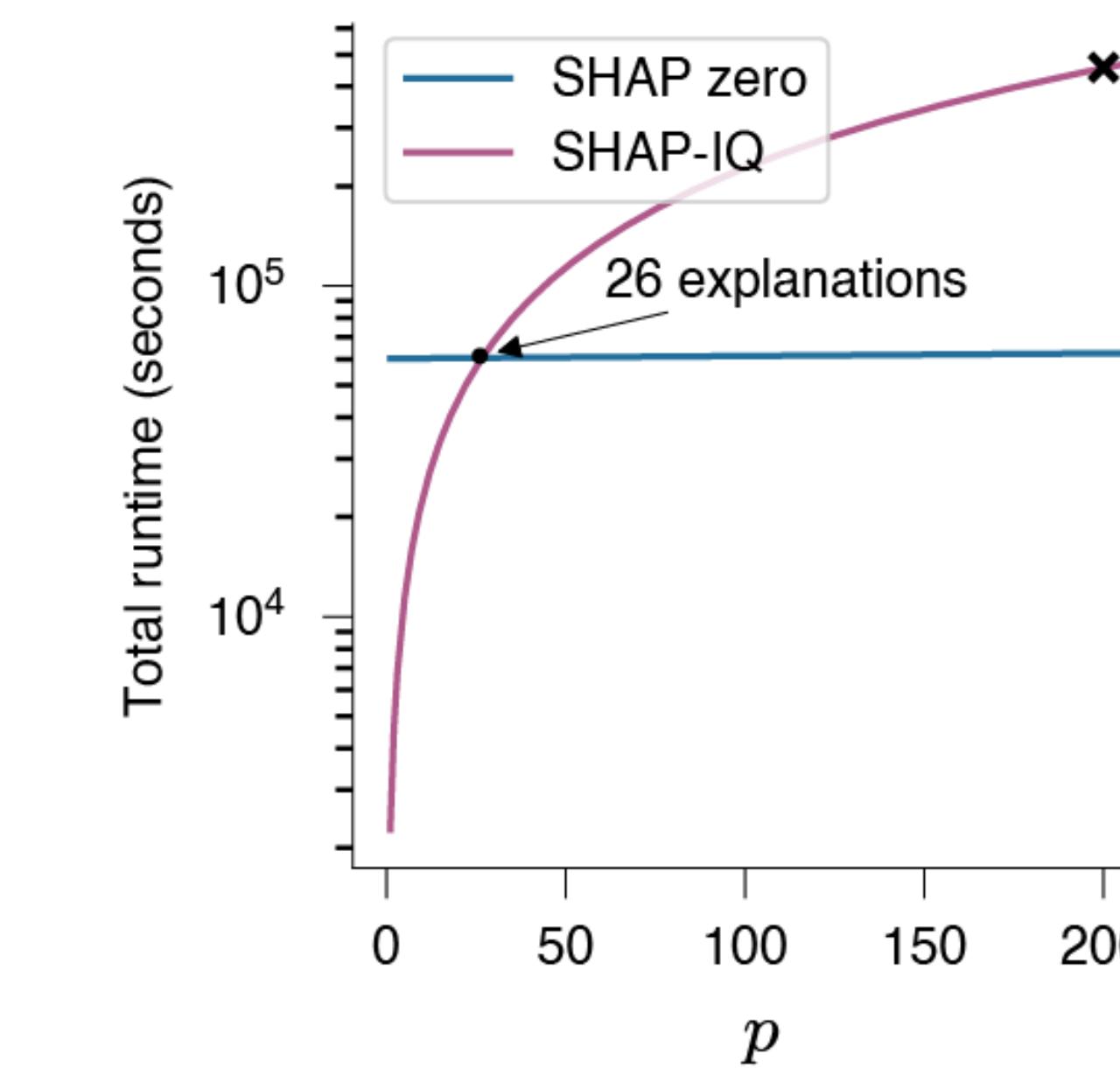
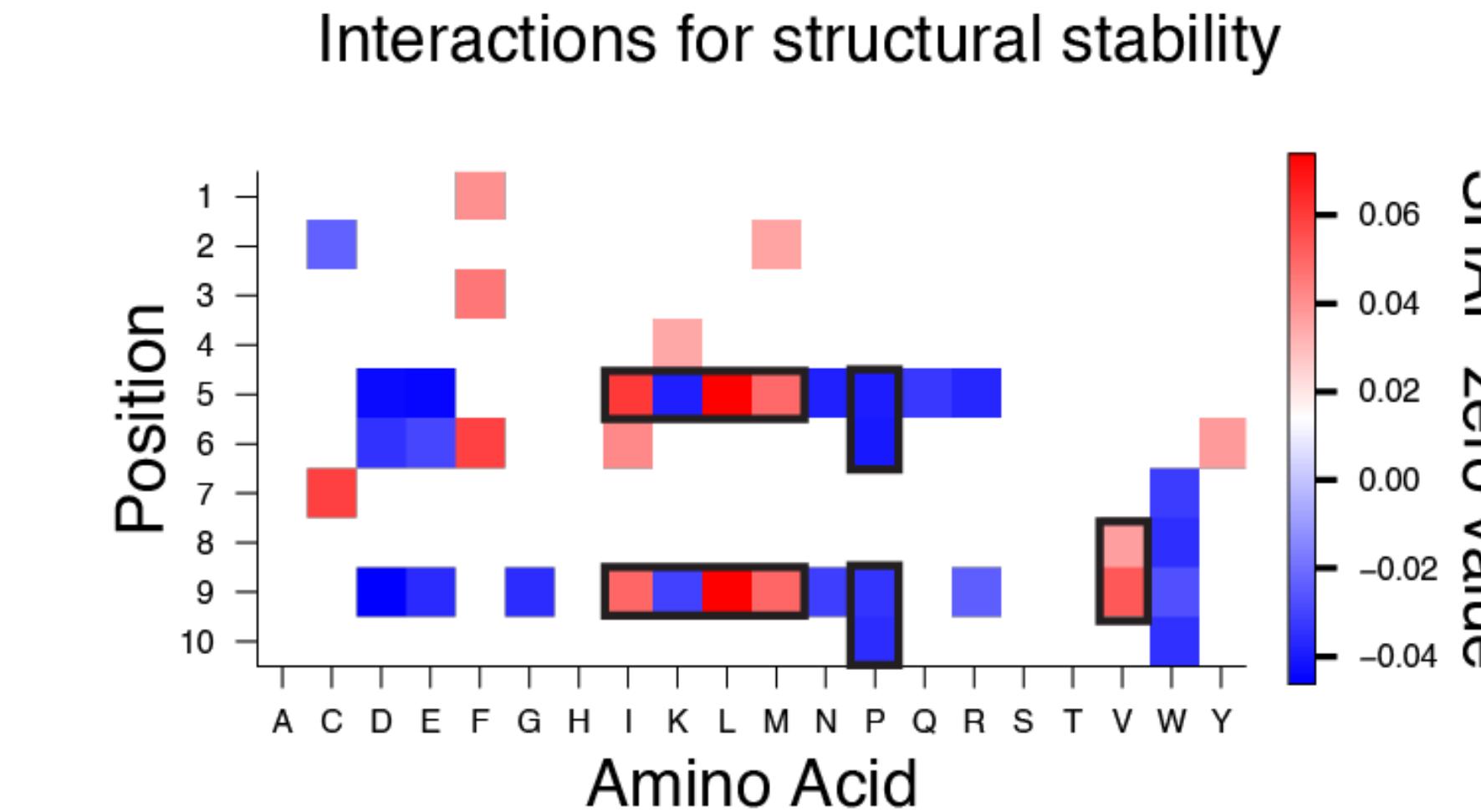
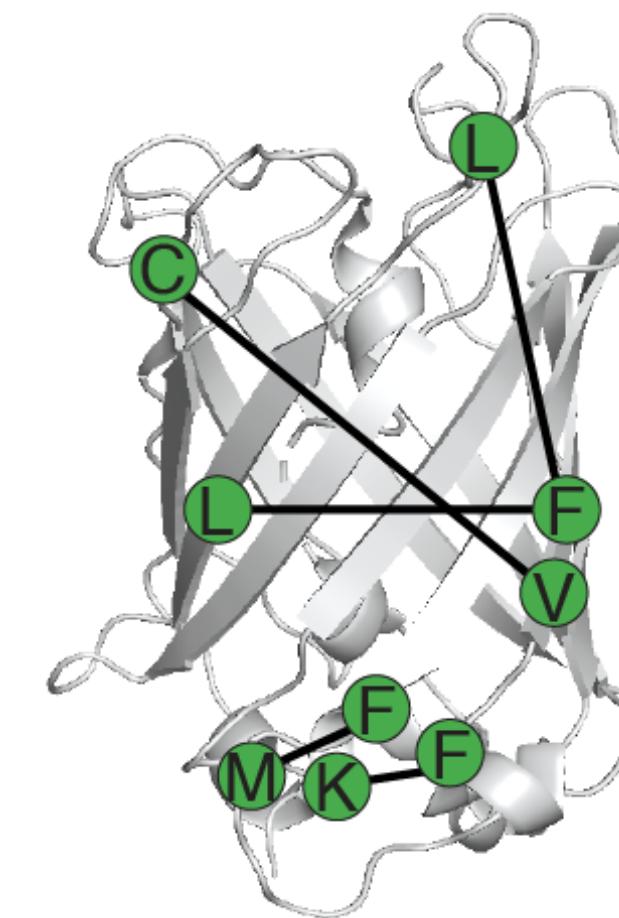
Interactions in guide RNA binding



# SHAP zero uncovers protein function interactions at scale



Interactions in protein function



# promoting open-source algorithms in biology

## SHAP zero: Explaining Biological Sequence Models

pypi package 0.0.4 license MIT status beta python 3.10 | 3.11 | 3.12 code style black  
last commit september

SHAP zero is a Python package that enables the amortized computation of Shapley values and interactions. It does this by paying a one-time cost to sketch the model's Fourier transform. After this one-time cost, SHAP zero enables **near-zero marginal cost** for future query sequences by mapping the Fourier transform to Shapley values and interactions.



pip install shapzero

# thank you!



Paper



Code