

# Base-resolution models of transcription factor binding reveal soft motif syntax

Žiga Avsec<sup>1</sup>, Melanie Weilert<sup>2</sup>, Avanti Shrikumar<sup>3</sup>, Sabrina Krueger<sup>2</sup>, Amr Alexandari<sup>3</sup>, Khyati Dalal<sup>2,5</sup>, Robin Fropf<sup>2</sup>, Charles McAnany<sup>2</sup>, Julien Gagneur<sup>1</sup>, Anshul Kundaje<sup>3,4\*</sup> and Julia Zeitlinger<sup>2,5\*</sup>

<sup>1</sup> Department of Informatics, Technical University of Munich, Garching, Germany

<sup>2</sup> Stowers Institute for Medical Research, Kansas City, MO, USA

<sup>3</sup> Department of Computer Science, Stanford University, Stanford, CA, USA

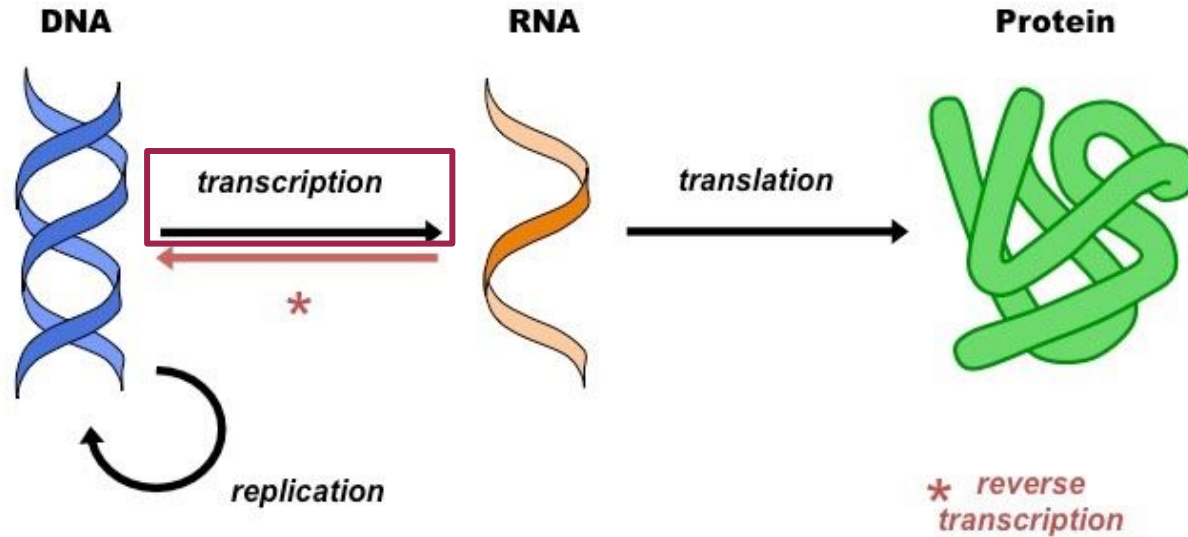
<sup>4</sup> Department of Genetics, Stanford University, Stanford, CA, USA

<sup>5</sup> The University of Kansas Medical Center, Kansas City, KS, USA

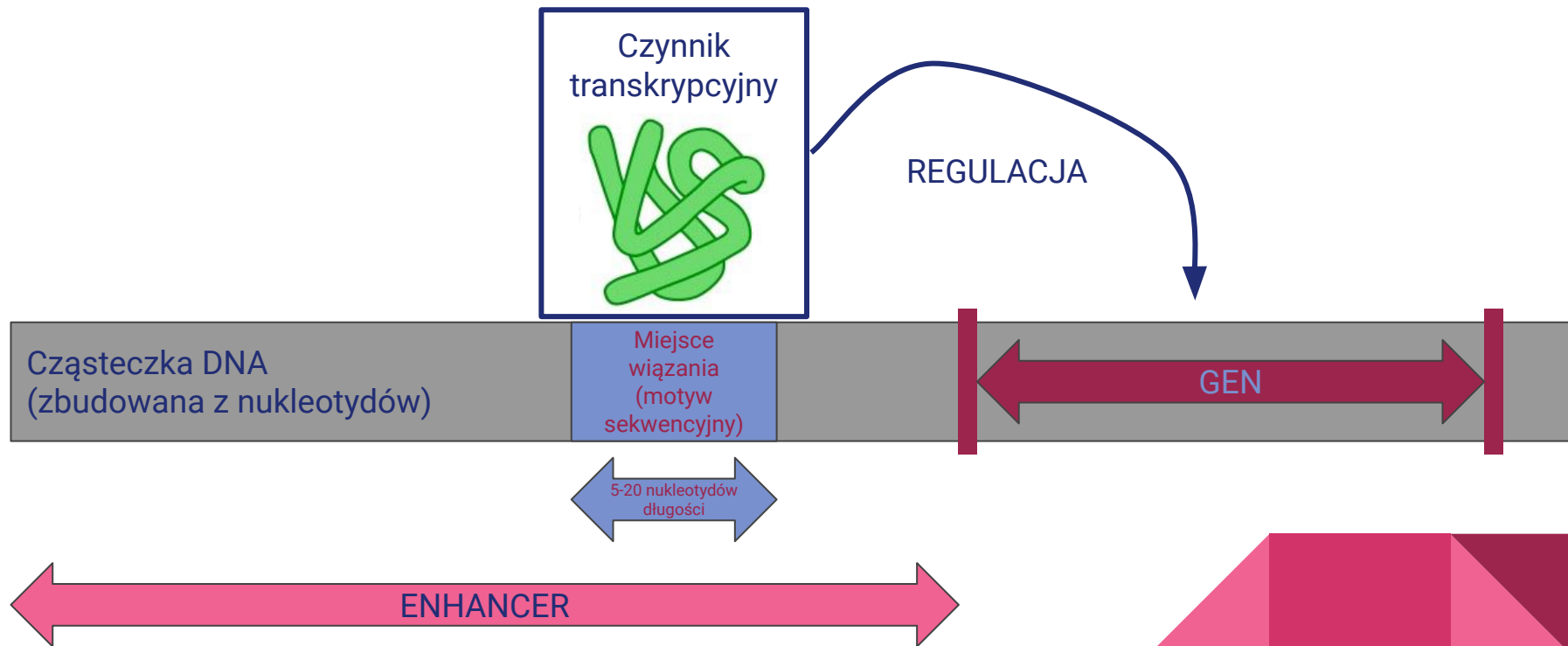
\* correspondence: [akundaje@stanford.edu](mailto:akundaje@stanford.edu), [jbz@stowers.org](mailto:jbz@stowers.org)

Autor prezentacji: Stanisław Antonowicz

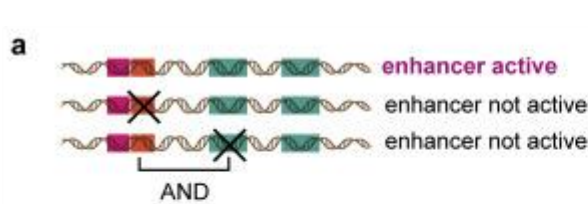
# Wstęp teoretyczny: centralny dogmat biologii molekularnej



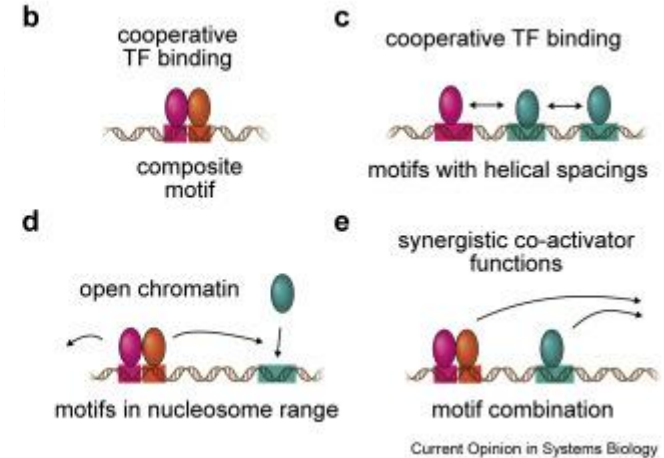
# Wstęp teoretyczny: czynniki transkrypcyjne



# Wstęp teoretyczny: składnia motywów



- Czynniki transkrypcyjne wiążąc się do genomu blisko siebie wchodzi z sobą w interakcje.
- Rozłożenie motywów sekwencyjnych nazywamy składnią.
  - “Silna” składnia to taka, gdzie motywy sekwencyjne oddziałują ze sobą w prosty sposób (np dwa motywy występują obok siebie). (b)
  - “Słaba” składnia to taka, gdzie ułożenie motywów na genomie jest bardziej skomplikowane bądź nieregularne. (c), (d), (e)



<https://doi.org/10.1016/j.coisb.2020.08.002>

# Wstęp teoretyczny: znaczenie czynników transkrypcyjnych

- Czynniki transkrypcyjne regulują ekspresję genów.
- Dzięki temu komórki mogą:
  - Reagować na zmiany w ich otoczeniu (organizmu bądź komórki).
  - Różnicować się w różne typy komórek.
  - Kontrolować cykl komórkowy.
- Niektóre czynniki transkrypcyjne są (częściowo) odpowiedzialne za niekontrolowane różnicowanie się i wzrost komórek rakowych.
  - Na przykład czynnik transkrypcyjny Nanog występuje w komórkach raka piersi.
    - Lu, X., Mazur, S., Lin, T. *et al.* The pluripotency factor nanog promotes breast cancer tumorigenesis and metastasis. *Oncogene* 33, 2655–2664 (2014). <https://doi.org/10.1038/onc.2013.209>

# Rozszyfrowanie znaczenia tytułu

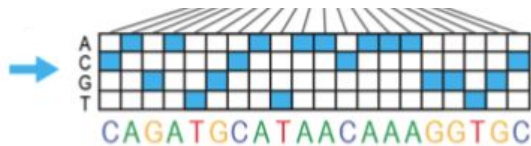
<b>Base-resolution</b> models	Modele z dokładnością co do nukleotydu
of <b>transcription factor binding</b>	wiązanie czynników transkrypcyjnych
reveal <b>soft motif syntax</b>	“słaba” składnia motywów

# Architektura modelu – wejście i wyjście

Wejście:

Sekwencje DNA w postaci macierzy wektorów “one-hot”.

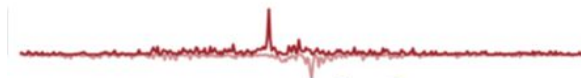
Input:  
~150,000  
sequences of 1 kb



Wyjście:

1. Profile częstotliwości wiązania czterech różnych czynników transkrypcyjnych w odcinku przed genem.
2. Ogólna liczba przyłączonych się cząsteczek danego czynnika transkrypcyjnego na danym odcinku.

Oct4



Sox2



Nanog



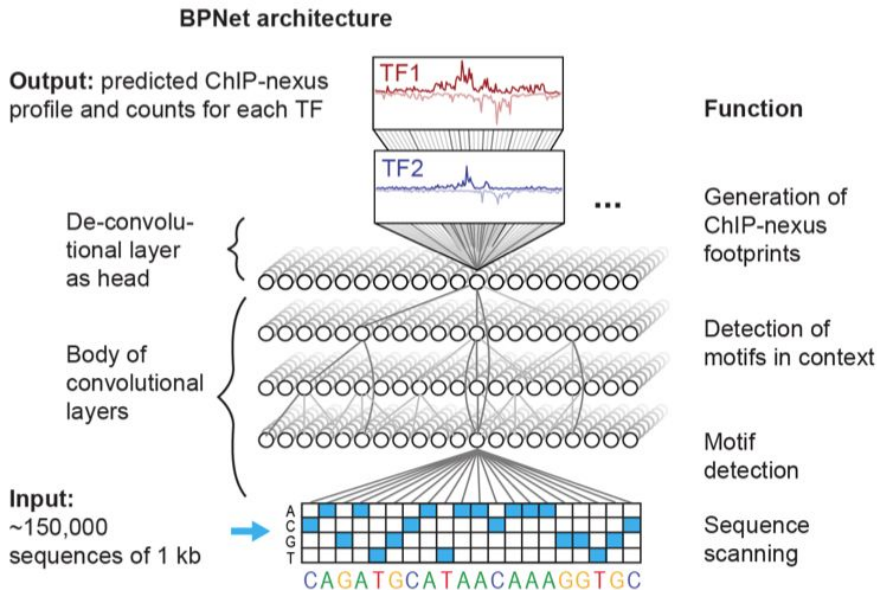
Klf4



# Architektura modelu

Kluczowe cechy modelu:

- Sieć konwolucyjna.
- Zastosowanie uczenia wielozadaniowego (*multitask learning*).
- Rezygnacja z poolingu.
- Kształt oraz łączna “masa” profilu częstotliwości są przewidywane oddzielnie. Pozwala na to niestandardowa funkcja straty.
- Używa pomocniczych danych o warunkach przeprowadzenia eksperymentu, żeby uniknąć biasu wynikającego z czynników losowych występujących w trakcie tworzenia zbioru danych.

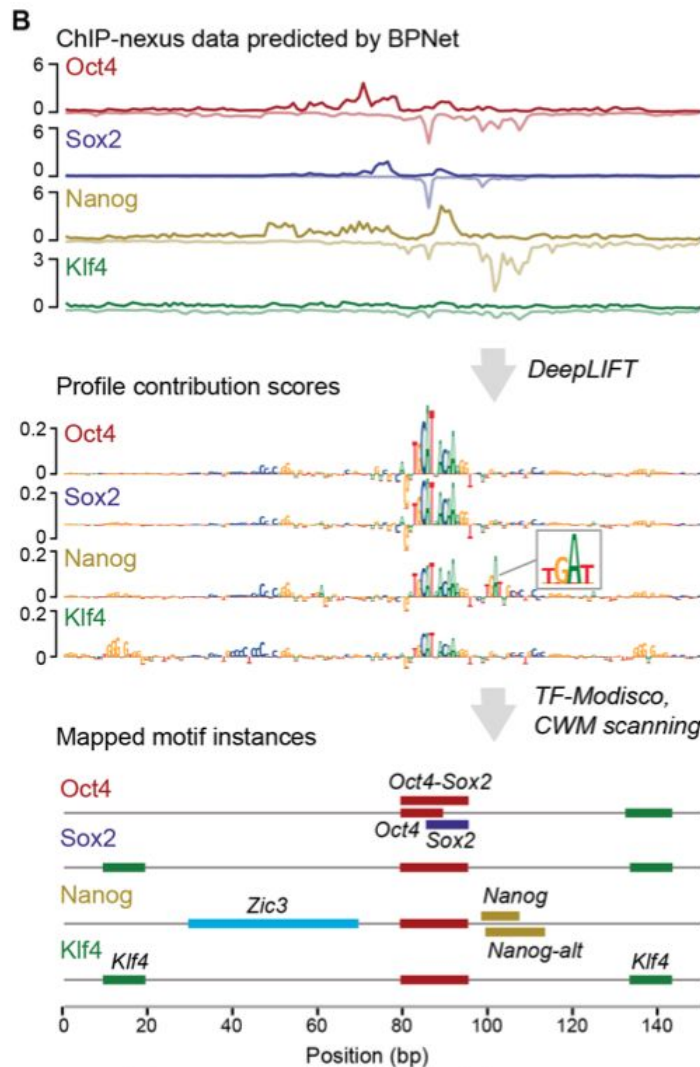


$$Loss = -\log p_{mult.}(\mathbf{k}^{obs} | \mathbf{p}^{pred}, n^{obs}) + \lambda(\log(1 + n^{obs}) - \log(1 + n^{pred}))^2$$



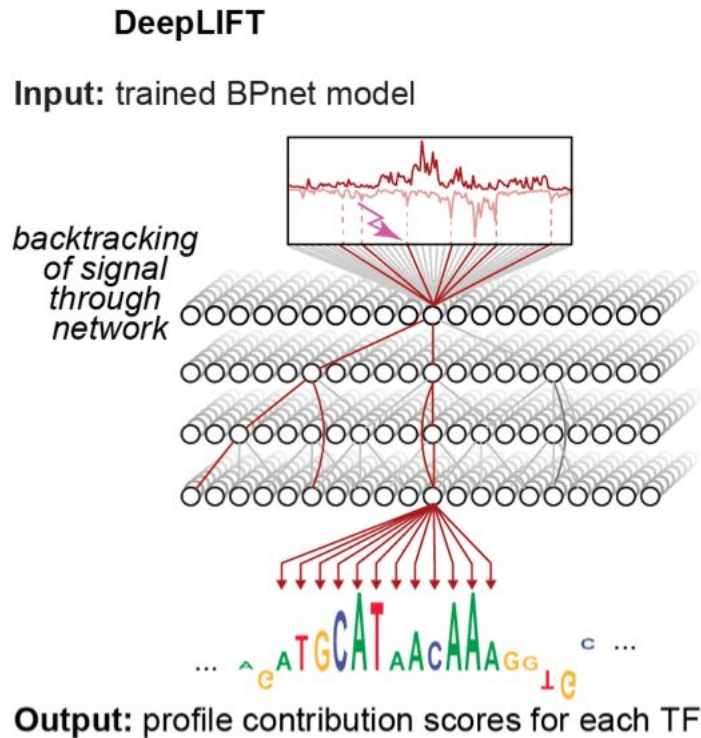
# Pipeline interpretacyjny

1. Predykcja profilu wiązania czynników transkrypcyjnych dla danego fragmentu DNA.
2. Profile kontrybucji dla każdego nukleotydu z sekwencji wejściowej.
3. Identyfikacja występowania konkretnych motywów w sekwencji wejściowej.



# Interpretacja modelu – DeepLIFT

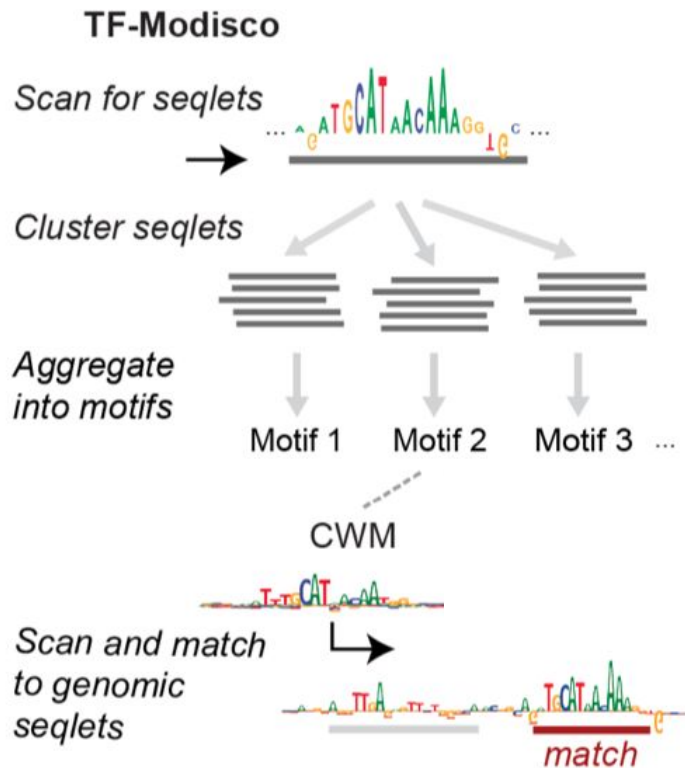
- Metoda gradientowa.
- Wyjaśniana jest różnica pomiędzy wyjściem sieci a wartością referencyjną (“neutralną”).
  - Dla MNIST może to być biały obraz, dla sekwencji DNA można użyć wielu permutacji sekwencji wejściowej i uśrednić wyniki.
  - Potrzebna jest wiedza domenowa, żeby zdefiniować dobrą referencję.
- Na drugim wykładzie ta metoda była wspomniana jako analogia do SHAPa.
- <https://github.com/kundajelab/deeplift>



# Interpretacja modelu – TF-MoDISco

(TF MOTif Discovery from Importance SCOrEs)

- Motywy sekwencyjne mogą się nieznacznie różnić od siebie, a mimo to być miejscami wiązań czynników transkrypcyjnych.
- TF-MoDISco pozwala na znalezienie konkretnych, znanych motywów sekwencyjnych w analizowanych sekwencjach.
  - Krótkie sekwencje z dużymi kontrybucjami są klastrowane.
  - Następnie dla każdego klastra jest obliczana reprezentacja CWM (contribution weight matrix), przedstawiająca kontrybucję każdego nukleotydu na każdej pozycji w klastrze.
  - Wejściowe sekwencje są skanowane z użyciem wyżej wspomnianej reprezentacji, żeby zidentyfikować wystąpienia konkretnych motywów w sekwencji.
- <https://github.com/kundajelab/tfmodisco>



# Wyniki – “silna” składania motywów

**A**

Single and composite motif's CWM

Structural basis

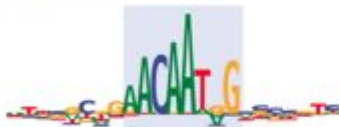
Oct4 (O1)



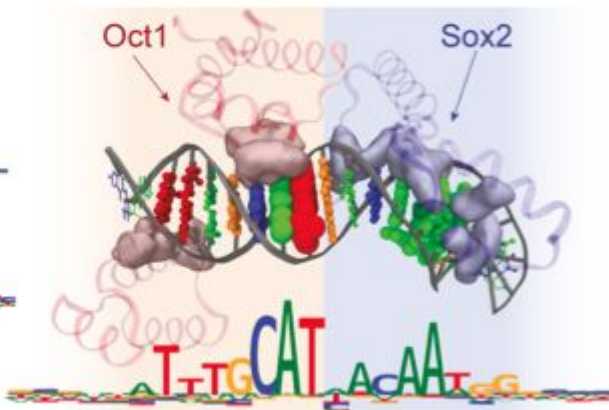
Oct4-Oct4 (O6)



Sox2 (S1)



Oct4-Sox2 (O0)



**b**

cooperative  
TF binding



composite  
motif

# Wyniki – “słaba” składania motywów

- Motywy *Oct4-Sox2* i *Sox2* wpływają na wiązanie czynnika transkrypcyjnego *Nanog*.
- Relacja jest jednokierunkowa – motyw *Sox2* wpływa na wiązanie *Nanog*, ale motywy *Nanog* nie wpływają na wiązanie *Sox2*.

C

All representative motifs

ID Motif CWM

O0 *Oct4-Sox2*



O1 *Oct4*



O6 *Oct4-Oct4\**



S1 *Sox2*



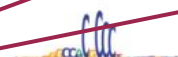
N1 *Nanog\*\*\**



N4 *Nanog-alt\*\*\**



K0 *Klf4*



K5 *Klf4-long\*\*\**



O5 *B-box\**



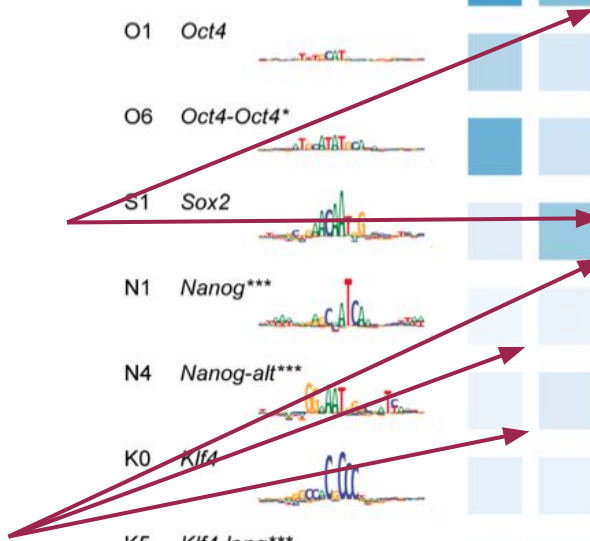
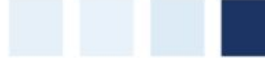
N2 *Zic3\**



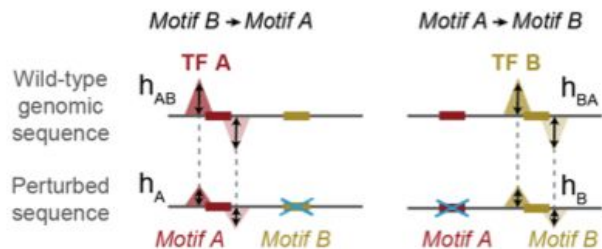
D

Average contribution score

Oct4 Sox2 Nanog Klf4

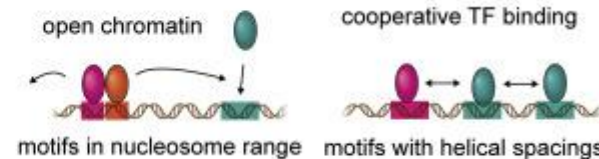
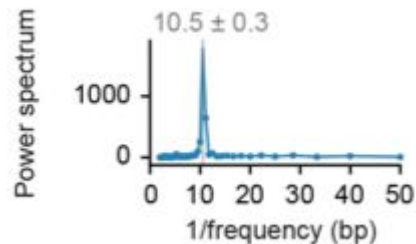
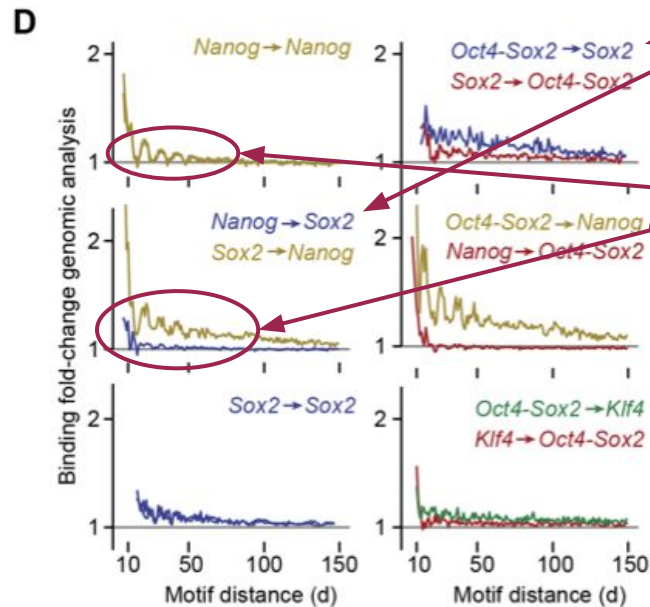


## B Genomic in silico mutagenesis interaction analysis



## Wyniki – “słaba” składania motywów.

- Niesymetryczność wpływu motywów sekwencyjnych na wiązanie czynników transkrypcyjnych.
- Okresowy charakter wiązania Nanog.



# Podsumowanie

- Konwolucyjne sieci neuronowe pozwalają na przewidywanie miejsc wiązania czynników transkrypcyjnych.
- Interpretacja modelu może nakierować badaczy na nowe fakty biologiczne.
- Wyniki pochodzące z modeli typu “black-box” mogą być “generatorem hipotez”, które można potem weryfikować eksperymentalnie.

