

Klasyfikacja obszarów funkcjonalnych

Wykład dla biotechnologów

Bartek Wilczyński
bartek@mimuw.edu.pl

12.12.2013

Modyfikacje histonów a chromatyna

The Histone Code

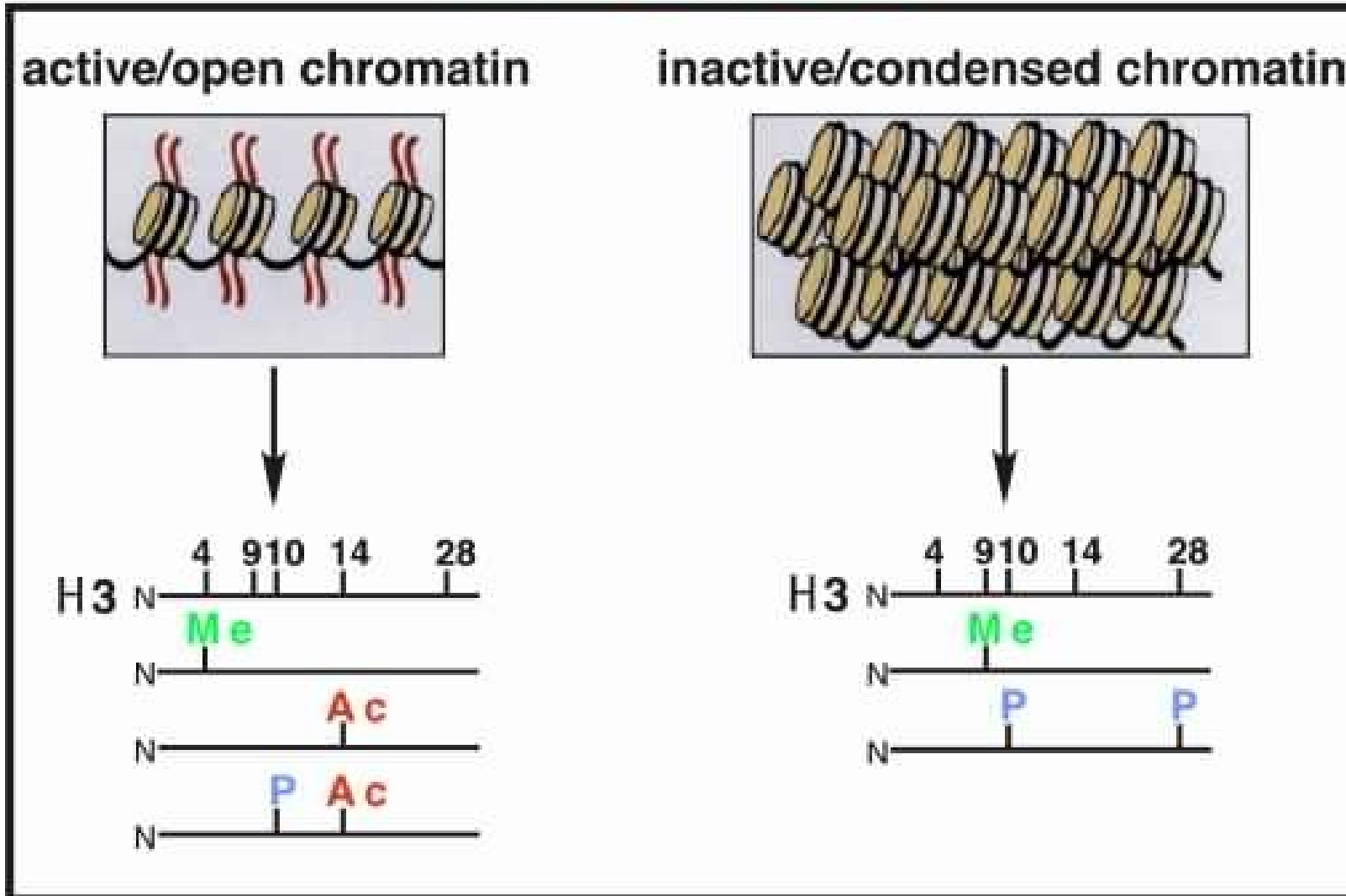
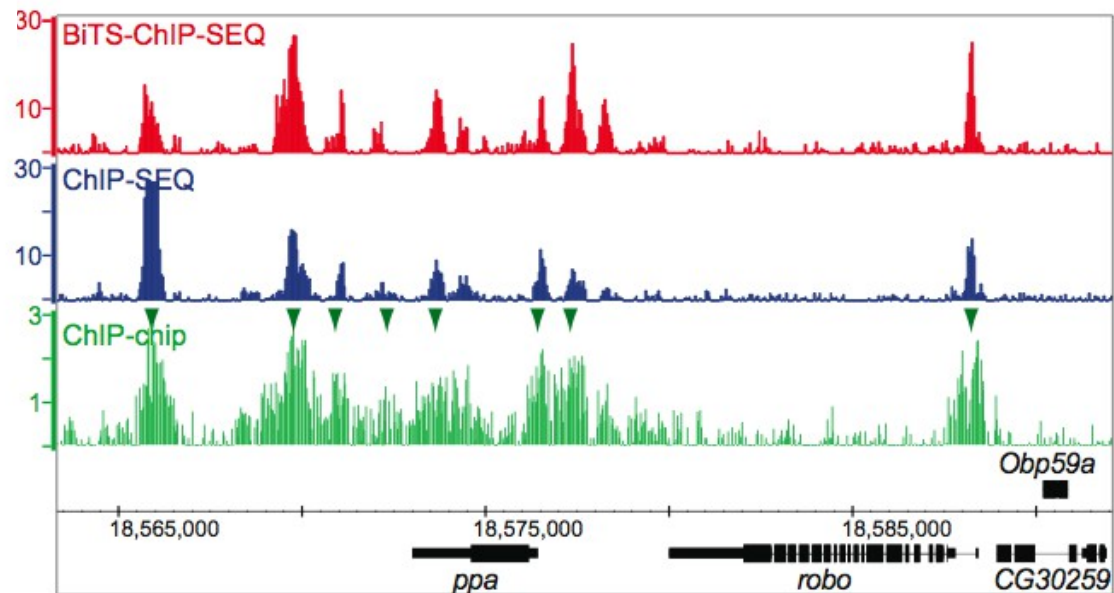
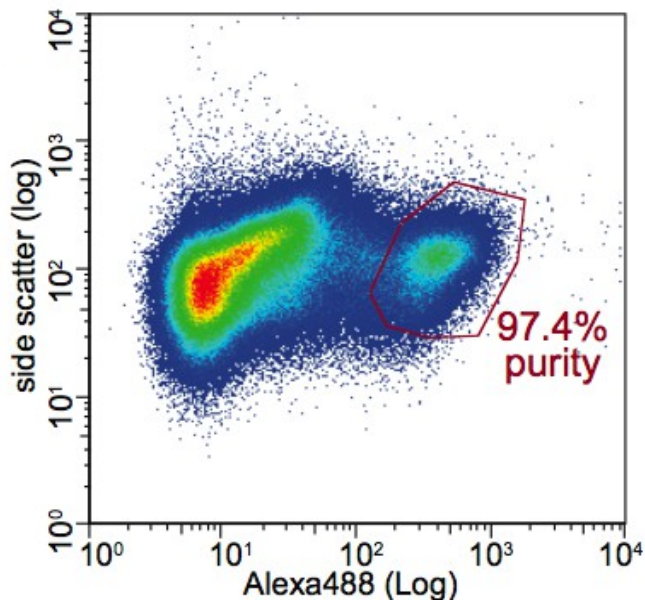
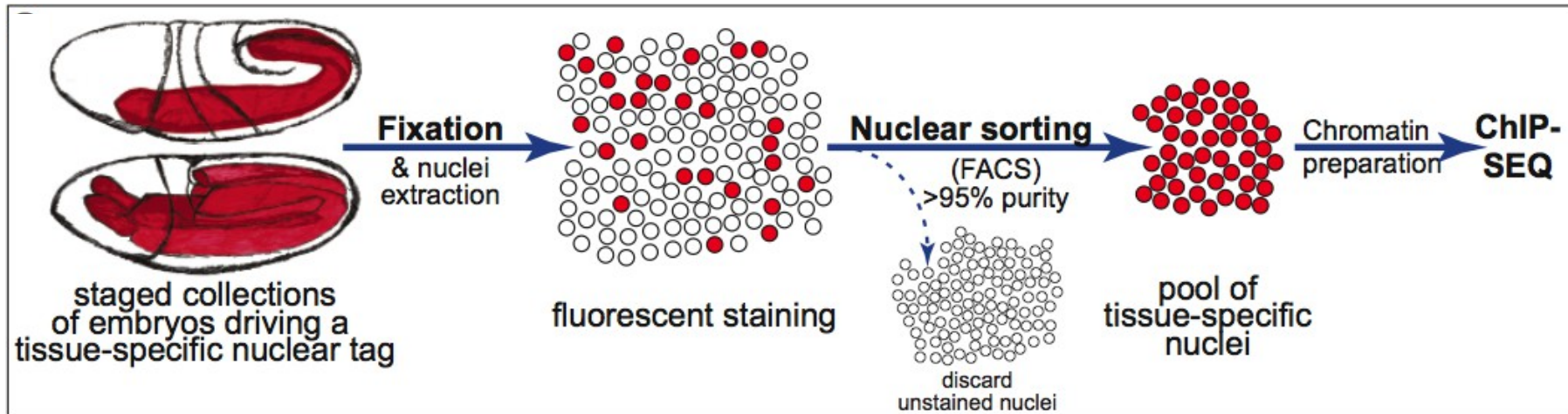


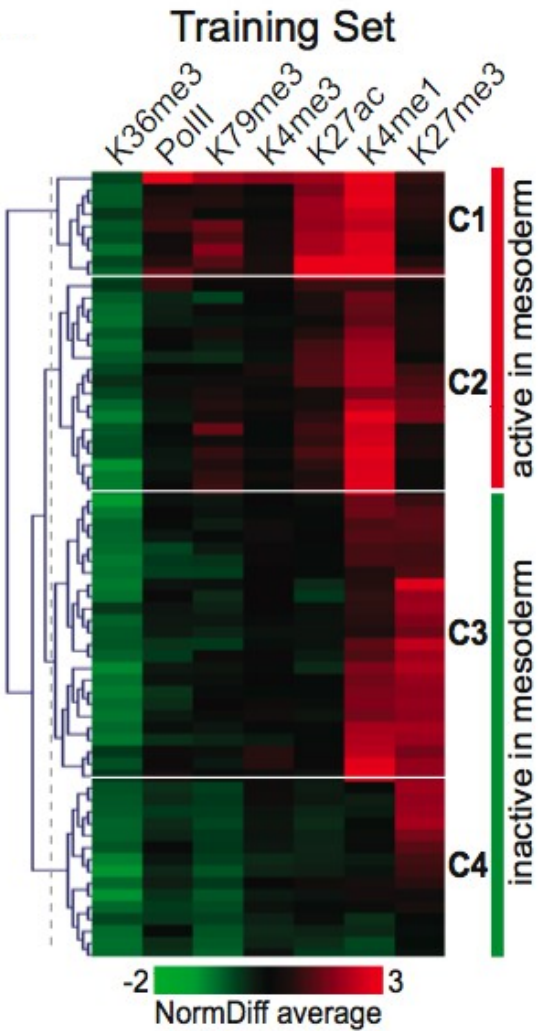
Figure 4

Modyfikacje specyficzne dla tkanki

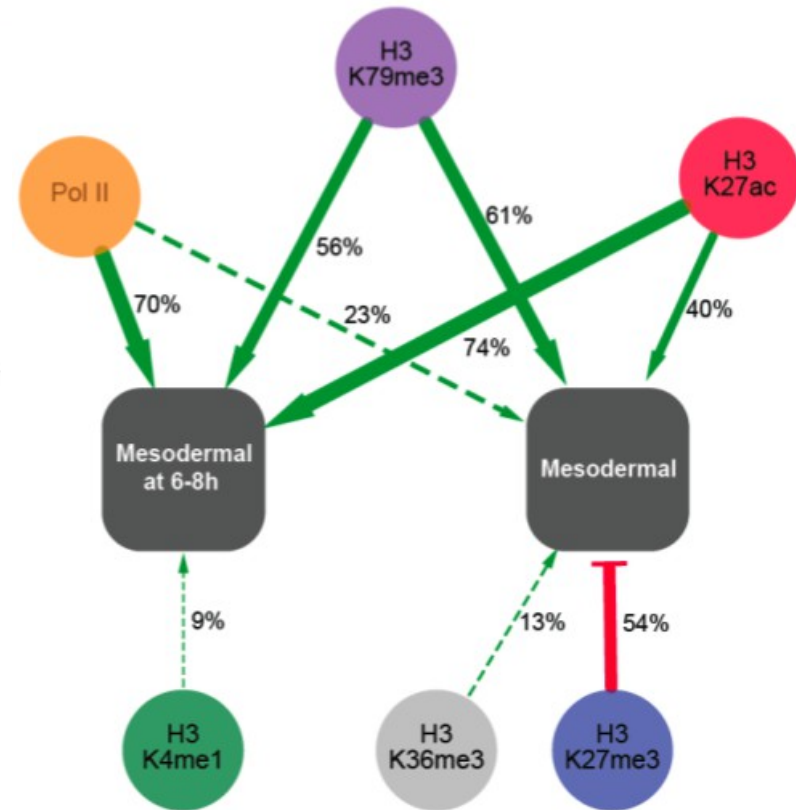
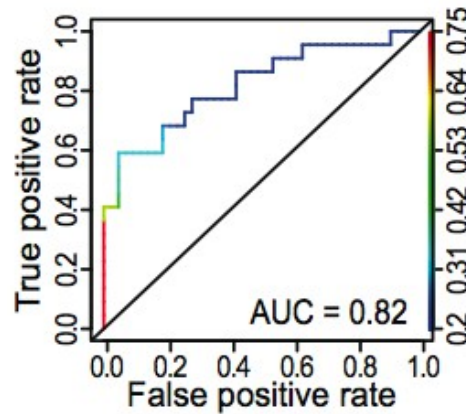


Bonn et al. Nat. Genet, 2012

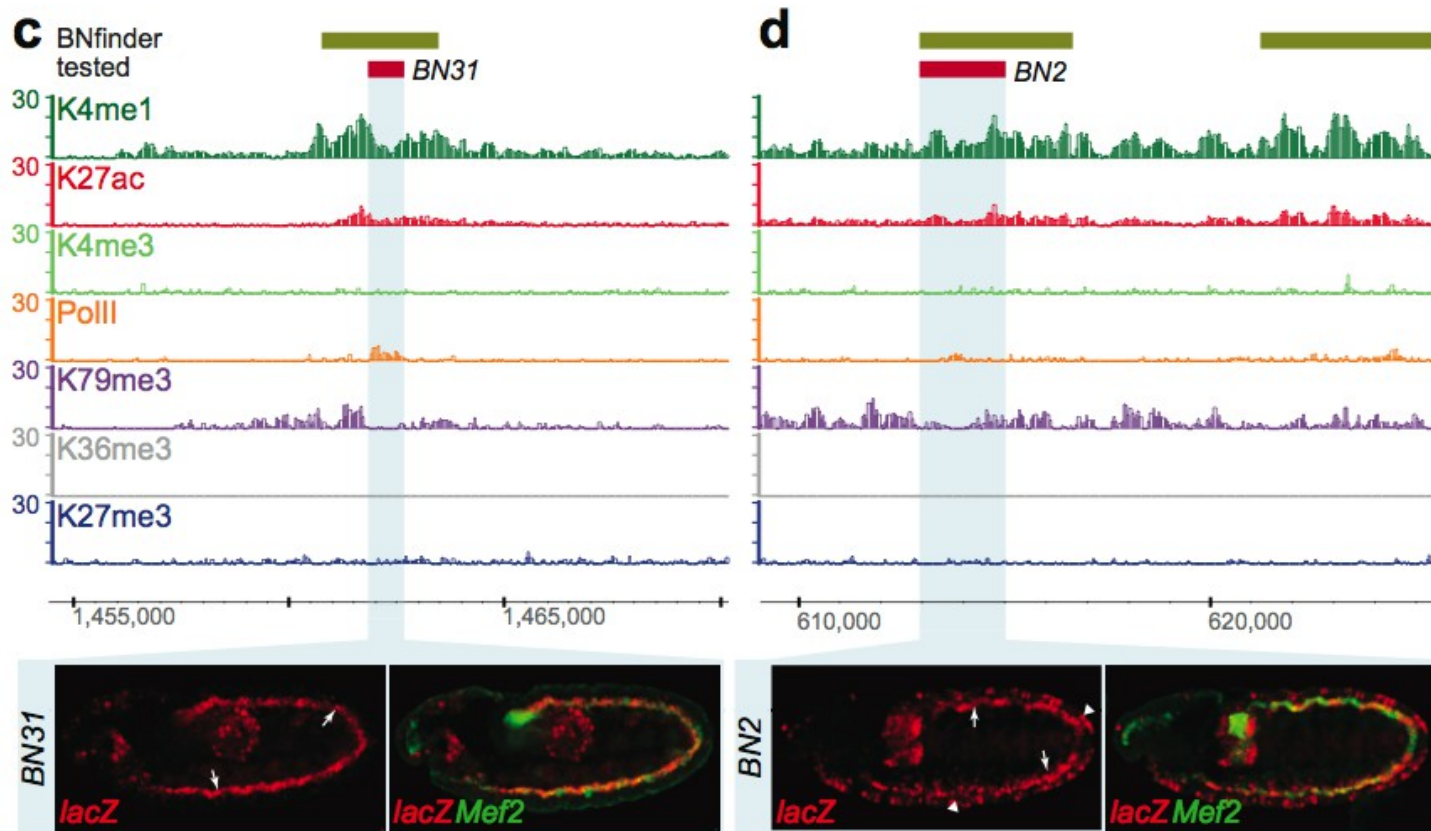
Stworzenie klasyfikatora



Bayesian Inference

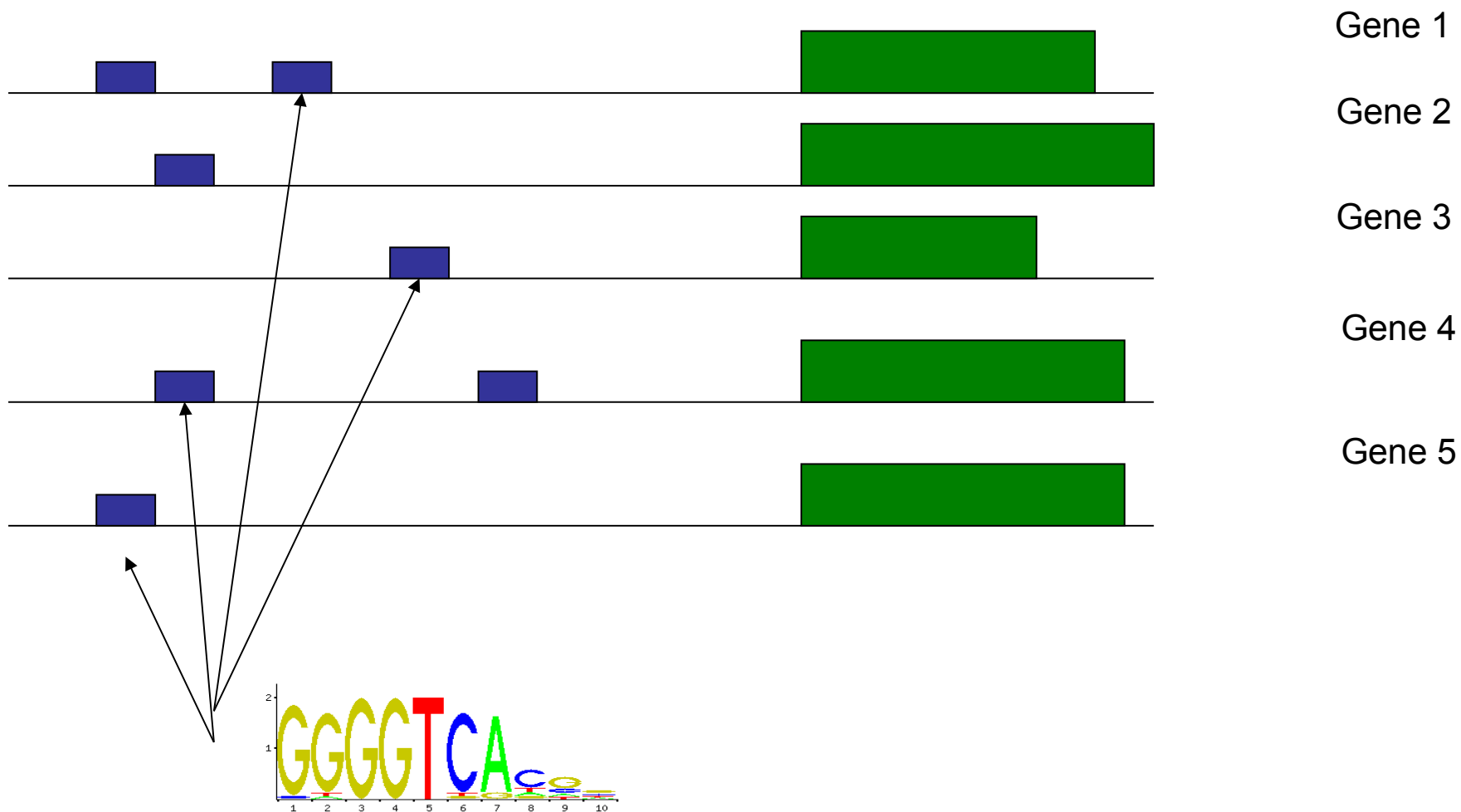


Weryfikacja eksperymentalna



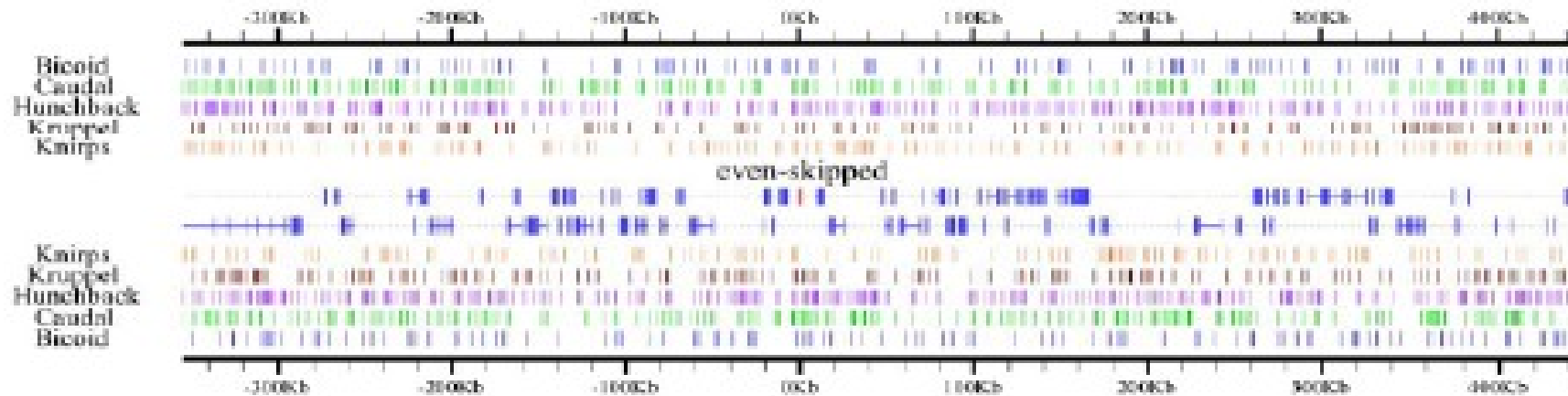
- 12 pozytywnych i 4 negatywne predykcje
- >90% prawidłowo! (1 pozytywna pomyłka)

Znajdowanie nowych motywów

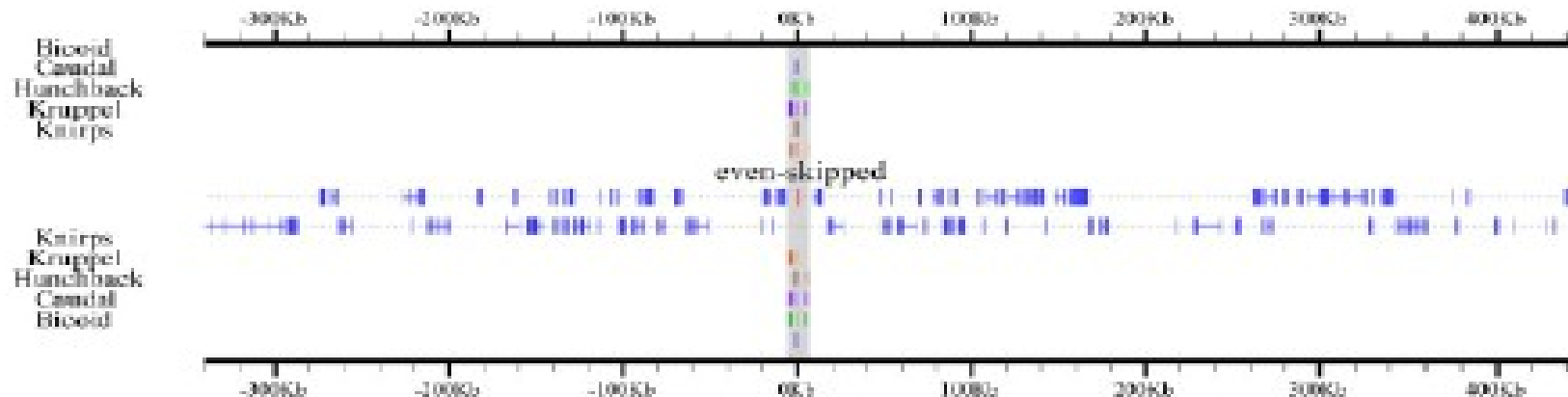


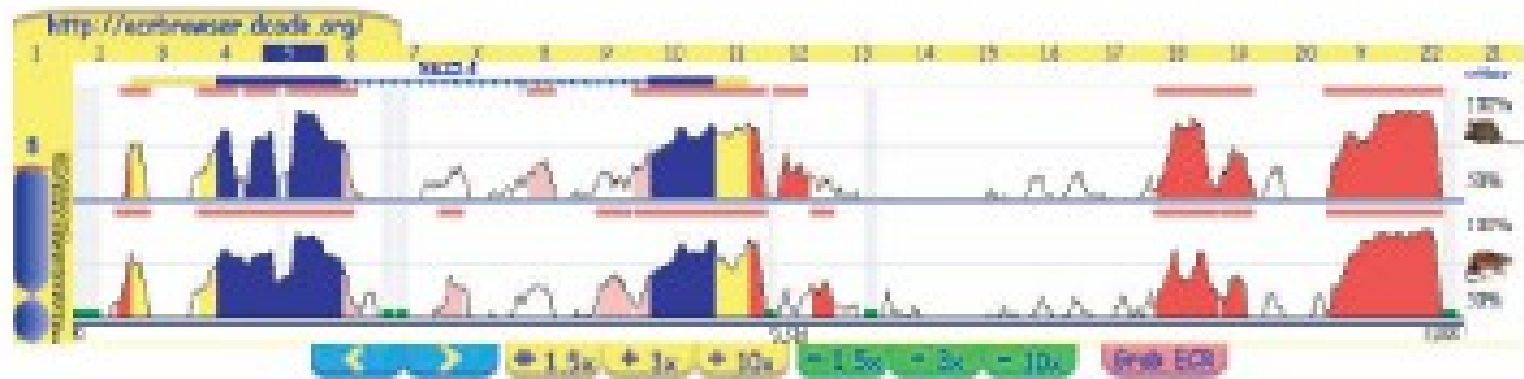
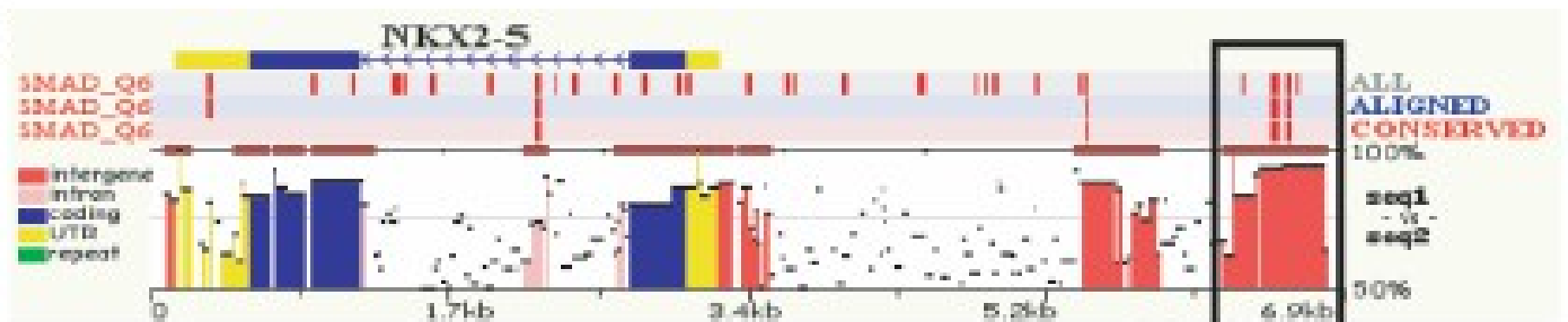
Szukanie Enhancerów przy pomocy motywów

(A) High stringency matches

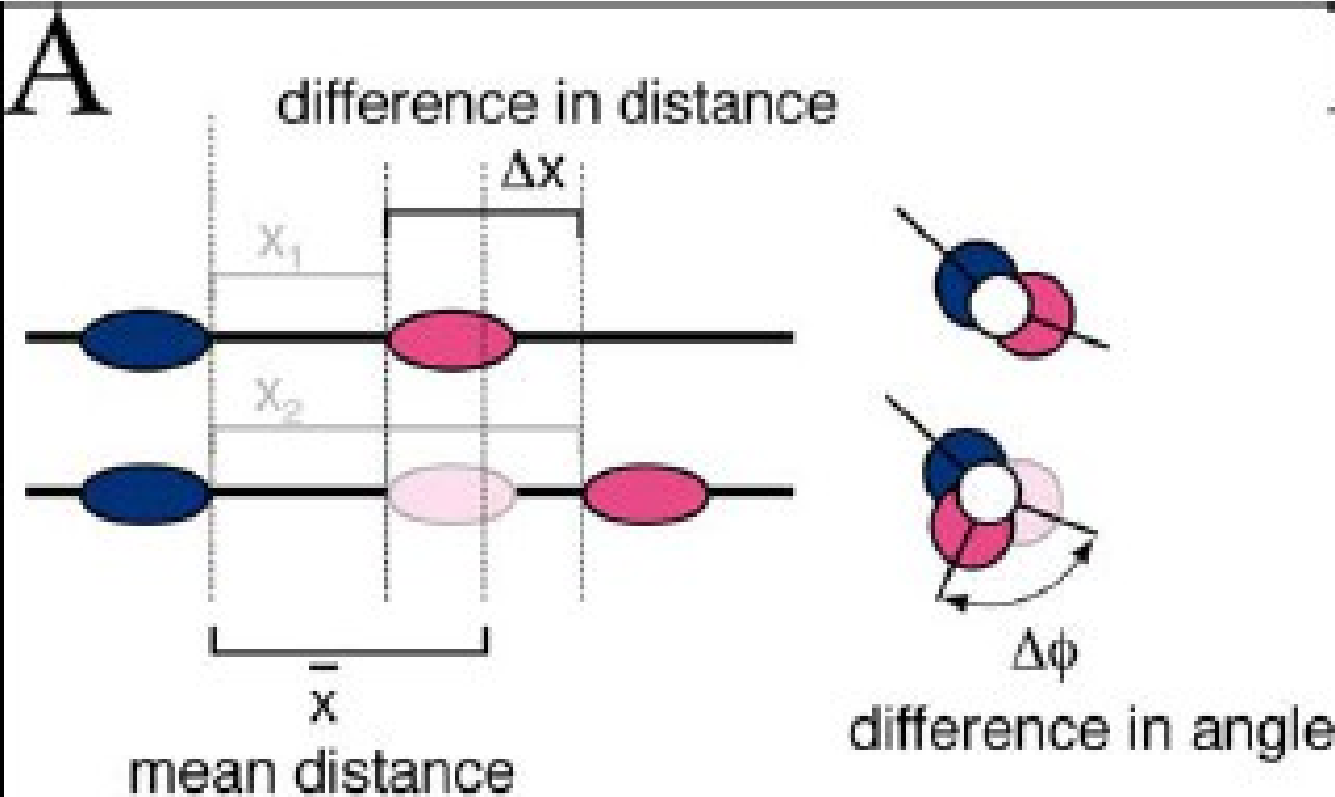


(B) High stringency matches and clustering filter



A**B****C**

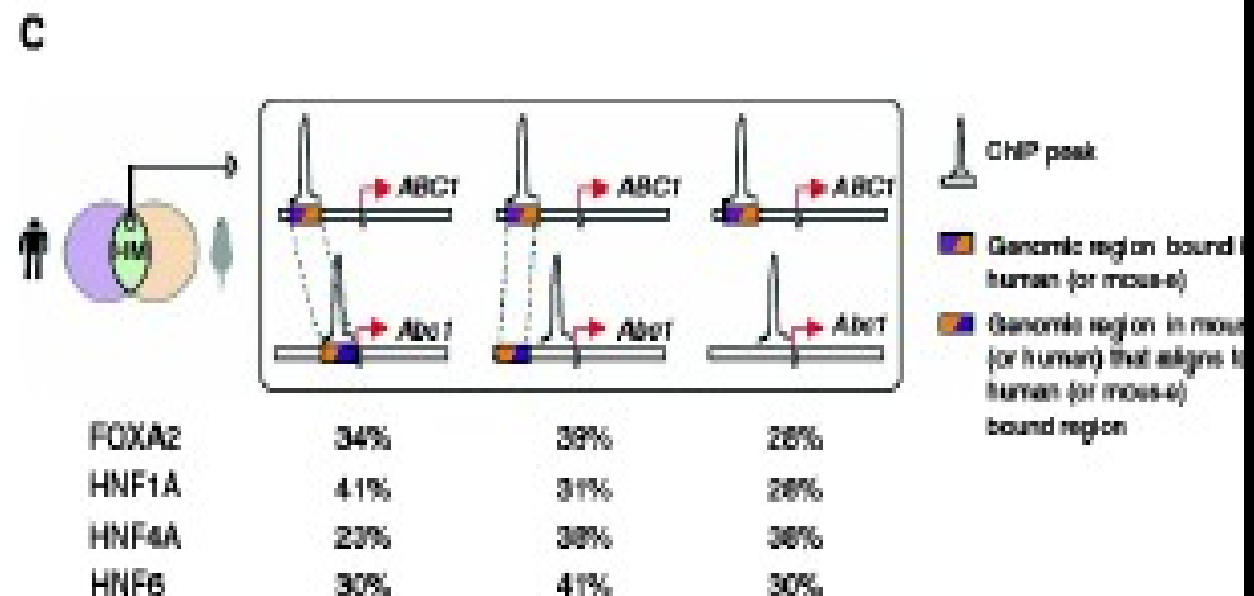
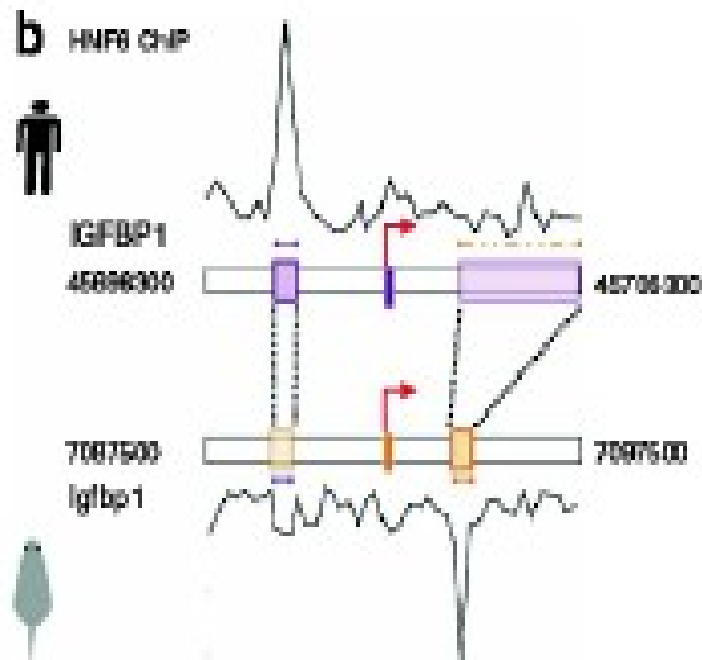
The EEL method (Halikias et al. Cell, 2006) aligns binding sites from promoter regions of homologous genes using a special function:



$$\text{Score} = \underbrace{\lambda \Delta G_T}_{\text{affinity}} - \underbrace{\mu \bar{x}}_{\text{clustering}} - \underbrace{\frac{\nu \Delta x^2 + \xi \Delta \phi^2}{2 \bar{x}}}_{\text{conservation}}$$

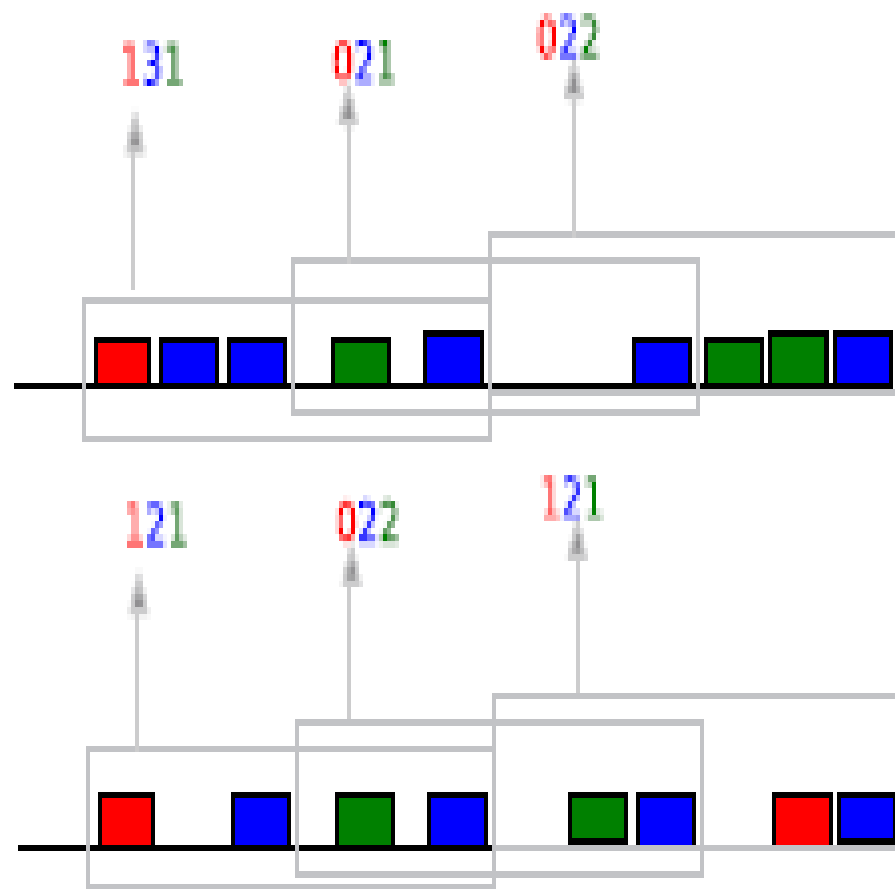
relative weights

Regulator	PFAM category	HS bound	MM bound	Intersection	Pvalue	HS binding sequence	MM binding sequence
FOXA2	Forkhead	151	574	68	1.0E-45		
HNF1A	POU-homeodomain	251	224	45	1.0E-29		
HNF4A	Nuclear receptor	1,251	654	387	1.0E-136		
HNF6	CUT-homeodomain	157	324	41	1.0E-27		



Odom et al. Nature Genetics, 2007, Similar results for Drosophila, Li et al, Genome Biol. 2007

$$S(P, P') = |P \cap P'| - \alpha \cdot |P \otimes P'| - \beta$$



$$P = p_1, p_2, p_3 = 131, 021, 022$$

$$P' = p'_1, p'_2, p'_3 = 121, 022, 121$$

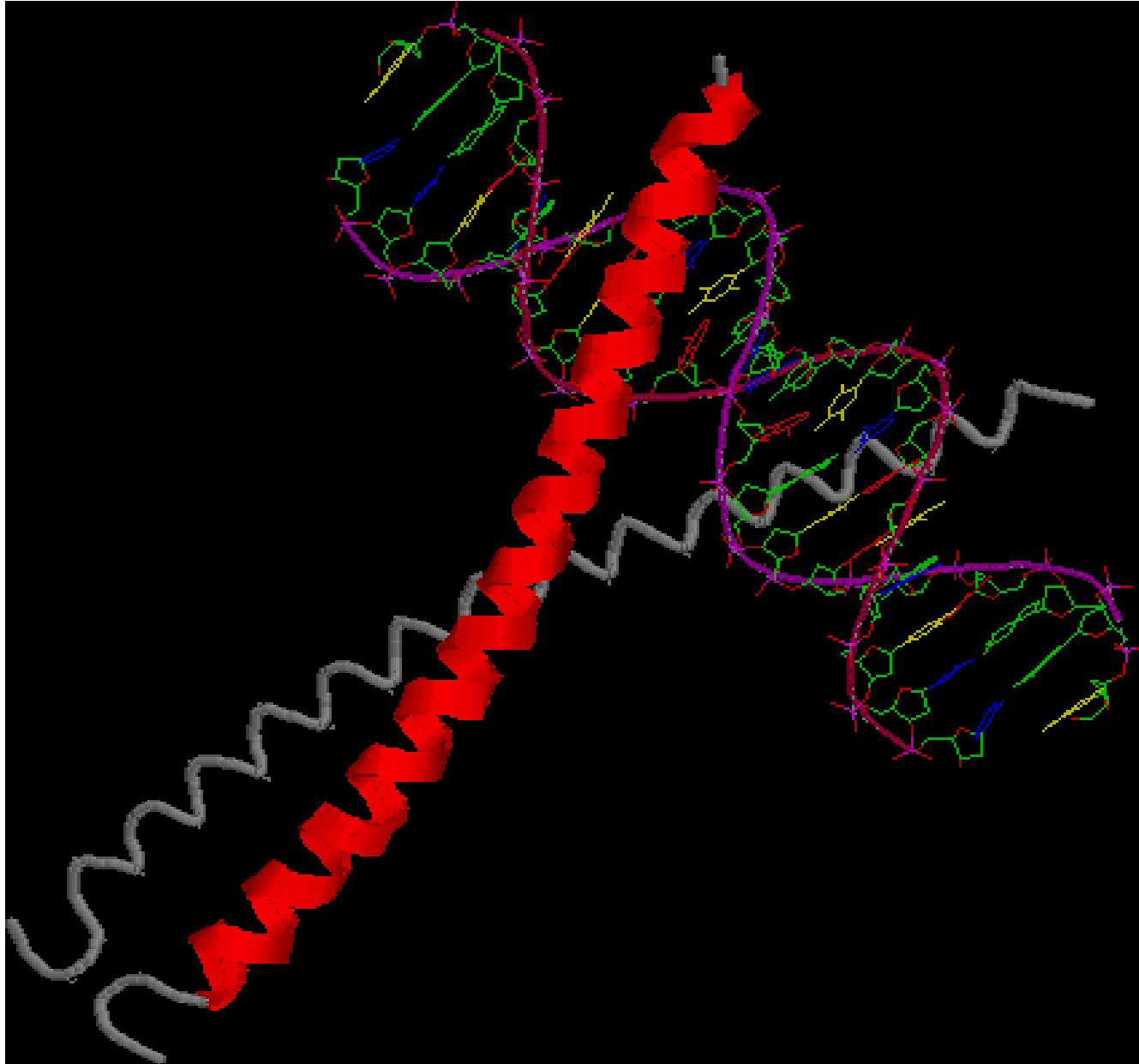
$$S(P, P') =$$

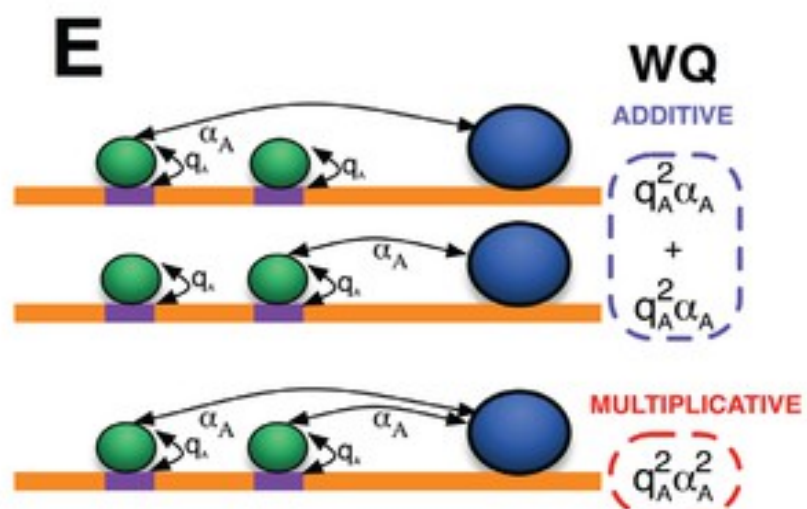
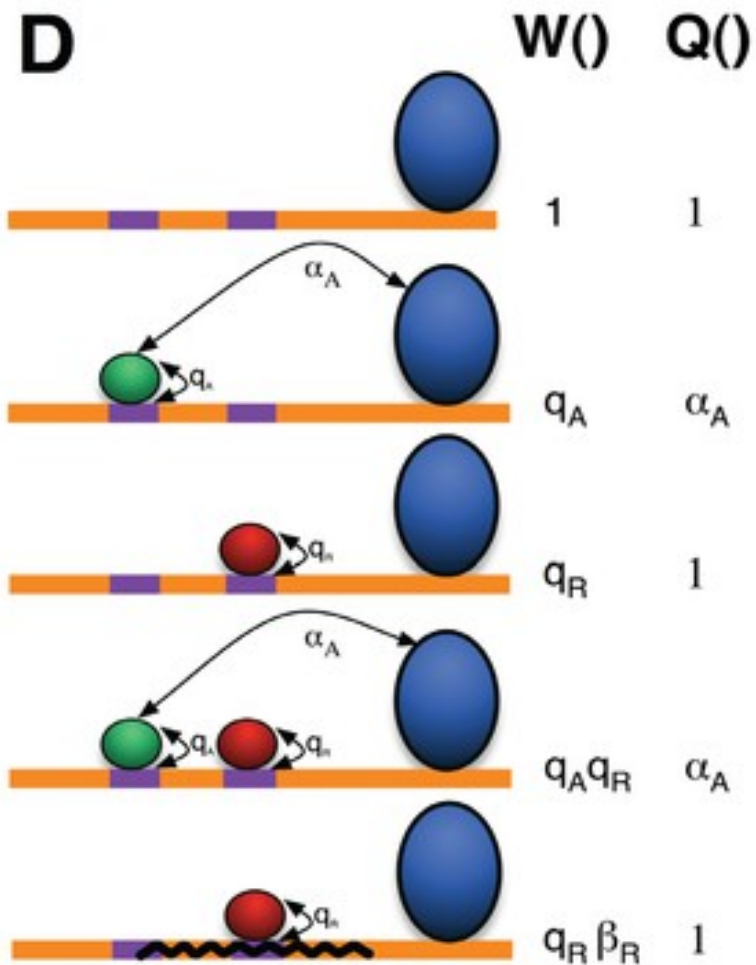
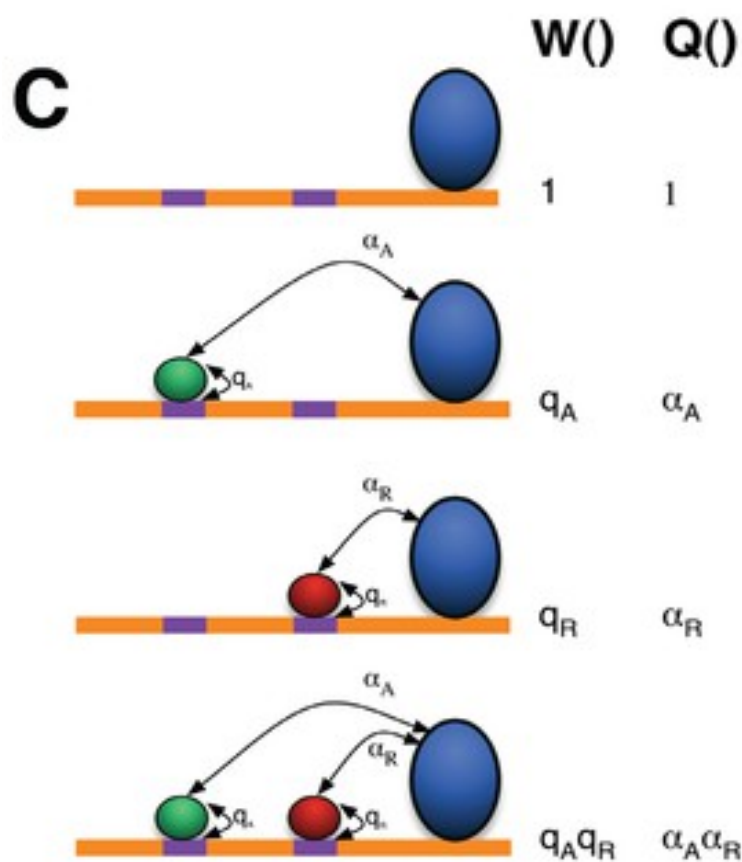
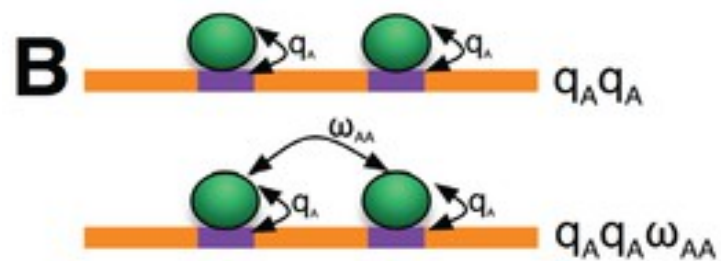
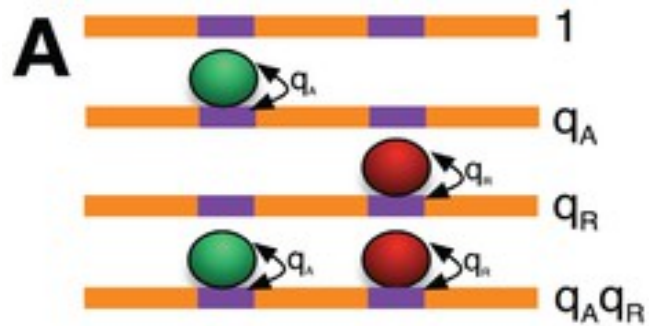
$$S(p_1, p'_1) + S(p_2, p'_2) + S(p_3, p'_3)$$

$$= (4 - \alpha - \beta) + (3 - \alpha - \beta) + (3 - 2\alpha - \beta)$$

$$= 9 - 4\alpha - 3\beta$$

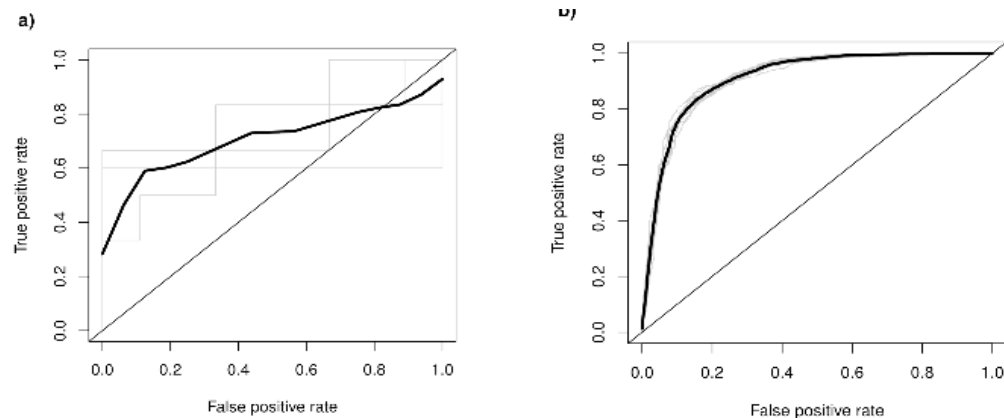
Wiązanie czynników transkrypcyjnych





We can add the sequence information

- First, we increase the size of the training dataset to use non-curated CRM predictions

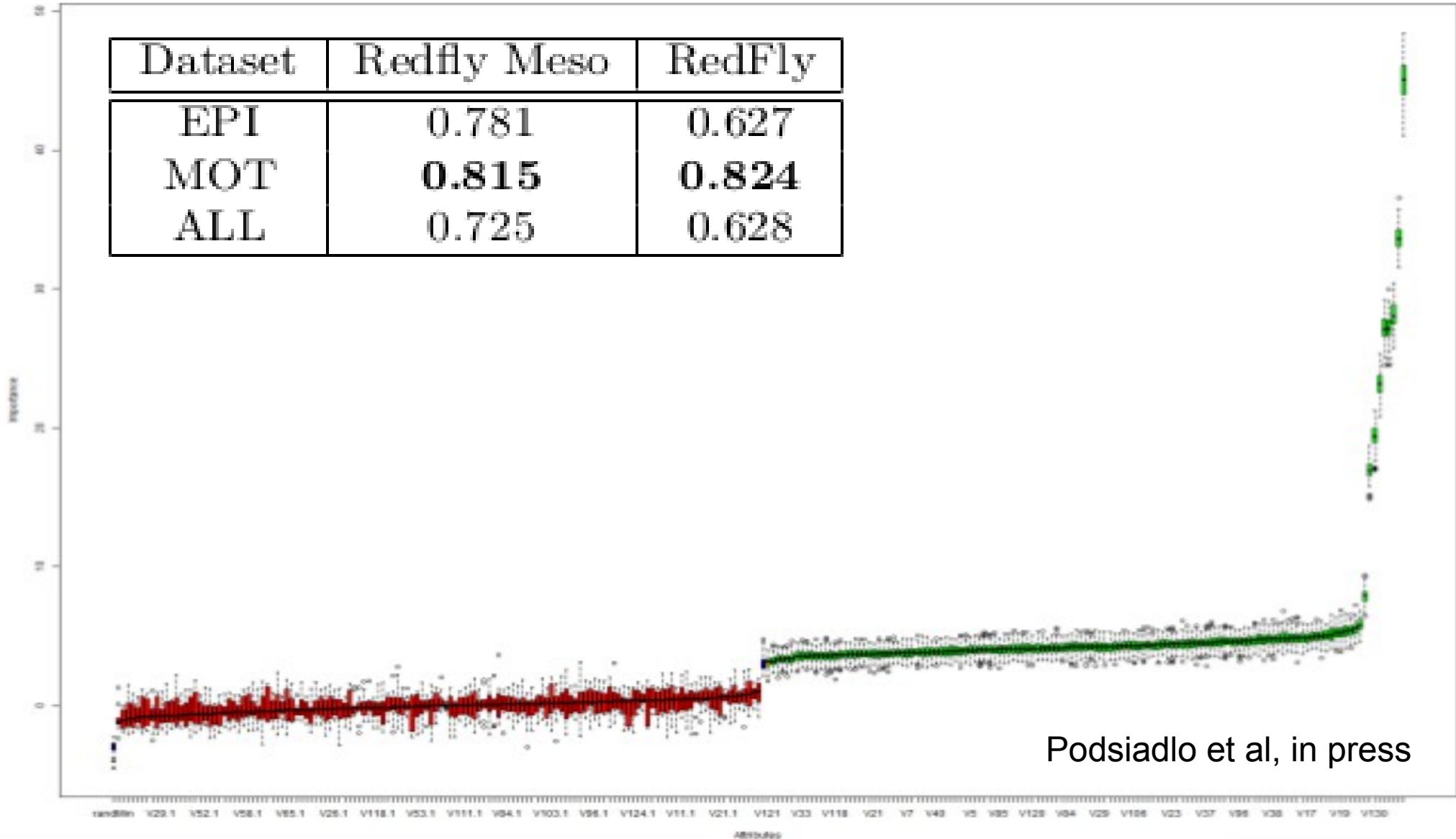


- Then, we add in data on motif affinity for JASPAR motifs

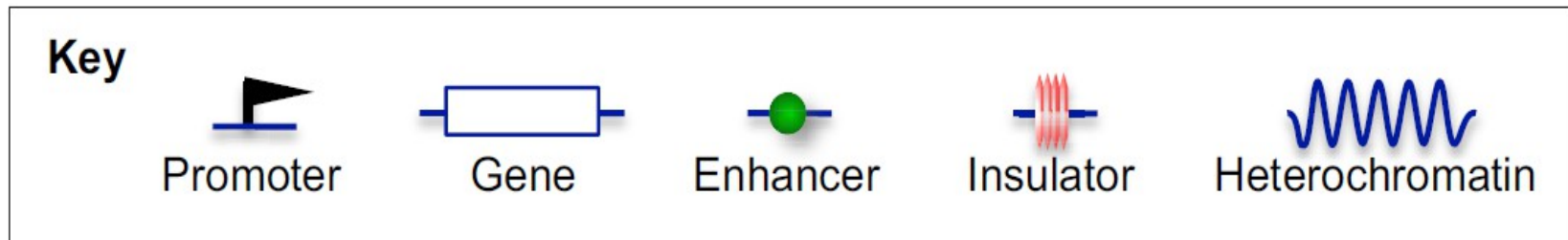
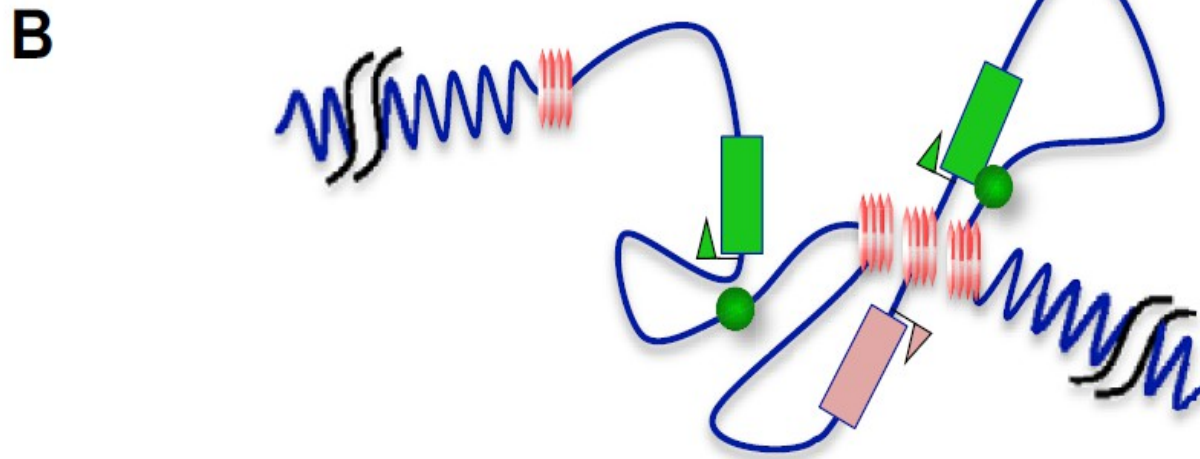
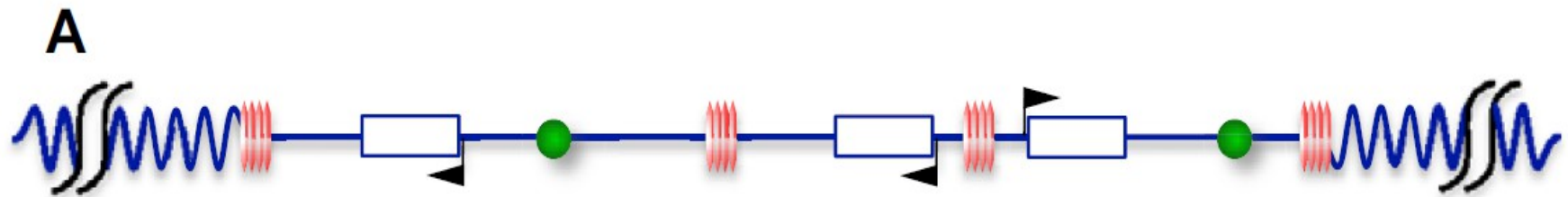
Dataset	BNFinder	SVM	RF
EPI	0.9	0.88	0.86
MOT	0.5	0.89	0.87
ALL	0.93	0.97	0.98

Importance of features is not consistent with generality

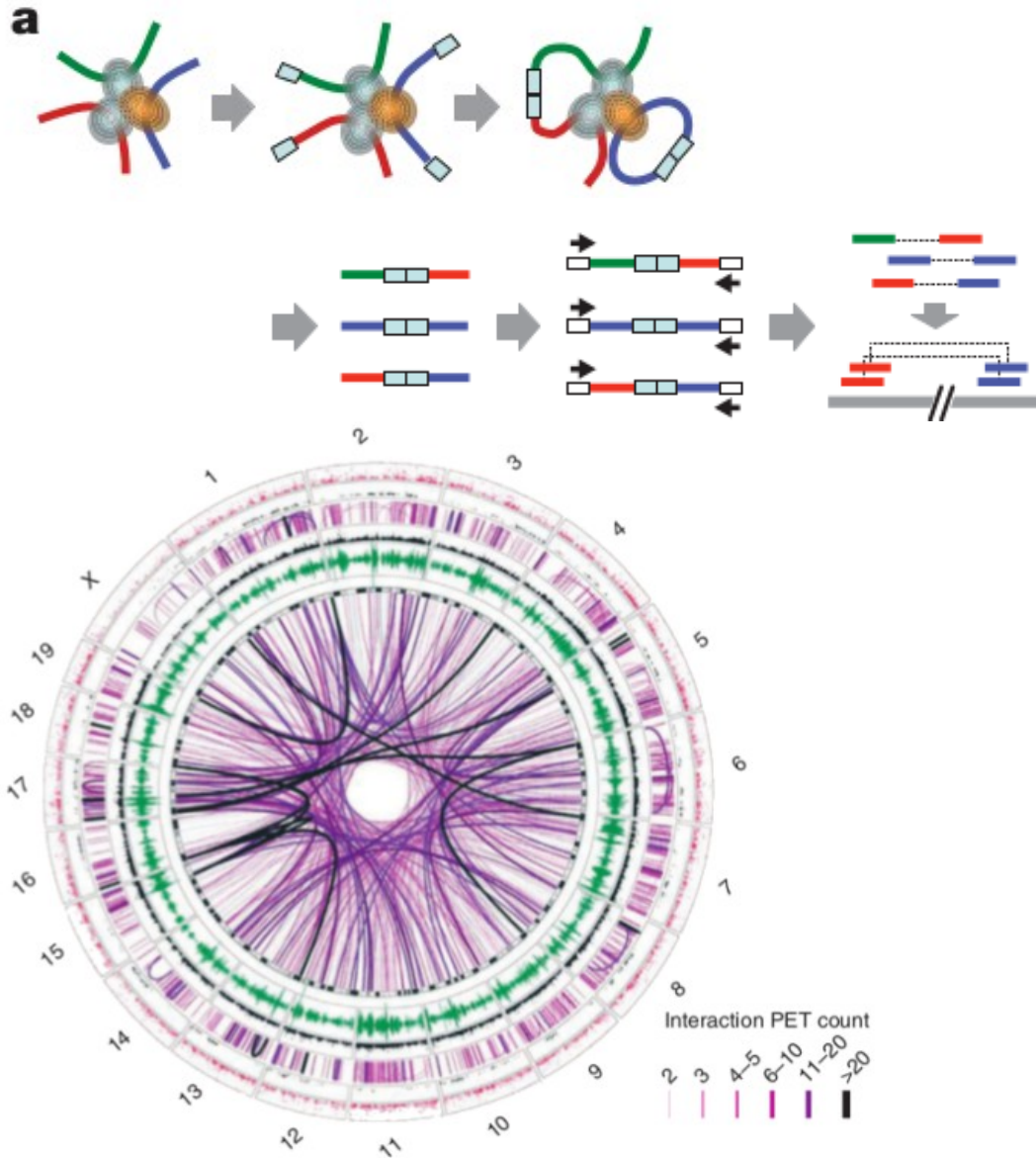
Dataset	Redfly Meso	RedFly
EPI	0.781	0.627
MOT	0.815	0.824
ALL	0.725	0.628



Insulator looping creates regulatory domains

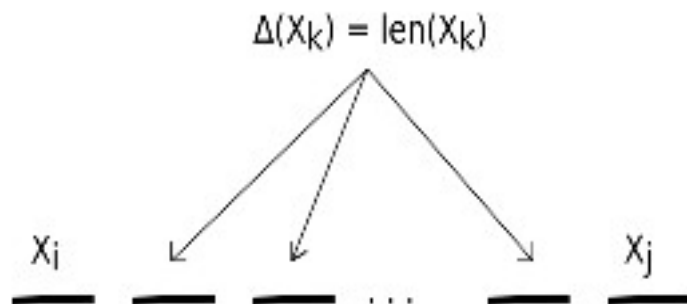


Enter chromosome interaction data



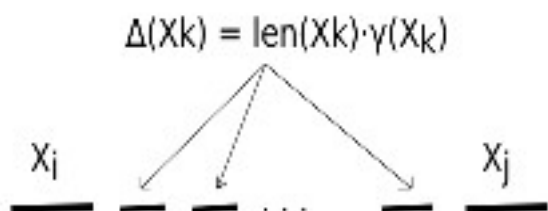
- Hi-C and ChIA-Pet protocols allow for measurement of chromatin interactions
- Getting this type of data is still difficult
- We can use computational methods to predict domains and interactions

Basic Model

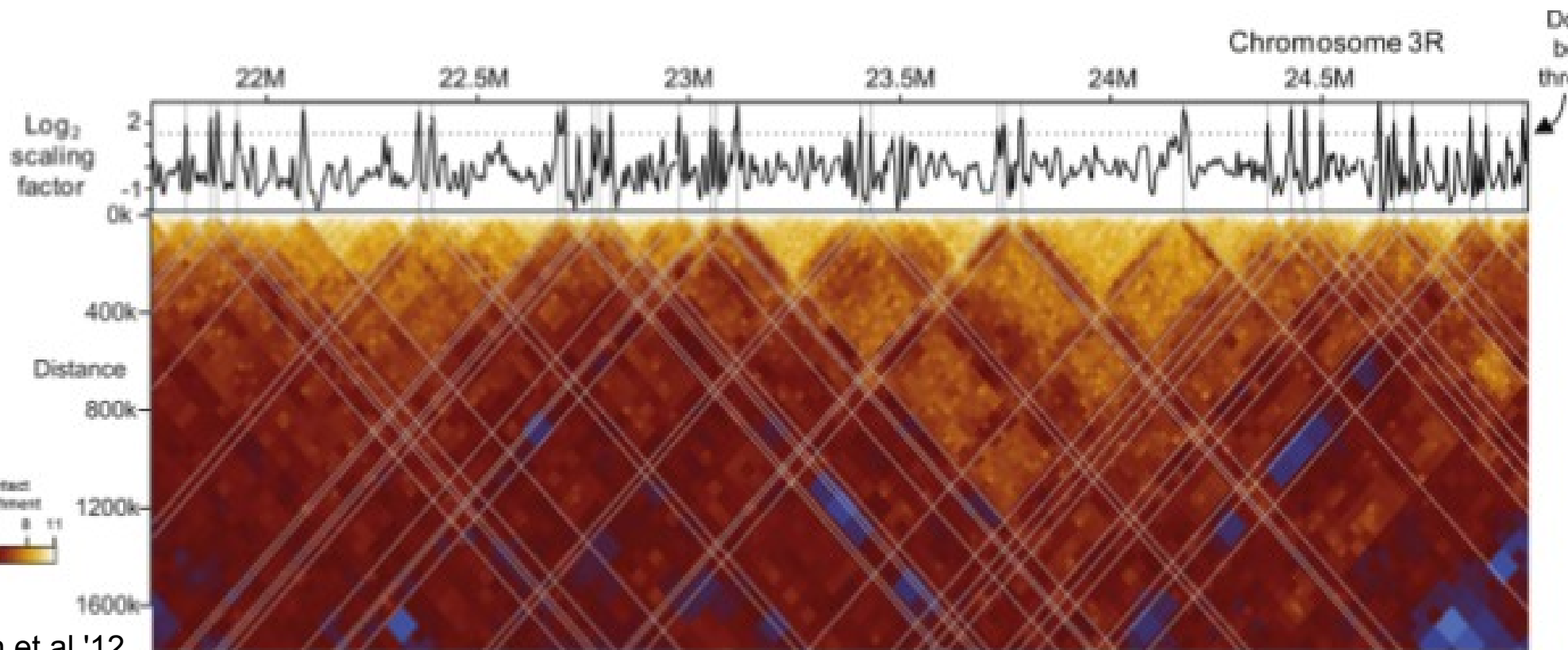


$$P(\text{contact}(X, Y)) = f\left(\sum_{i < k < j} \text{len}(X_k)\right)$$

Distance scaling model

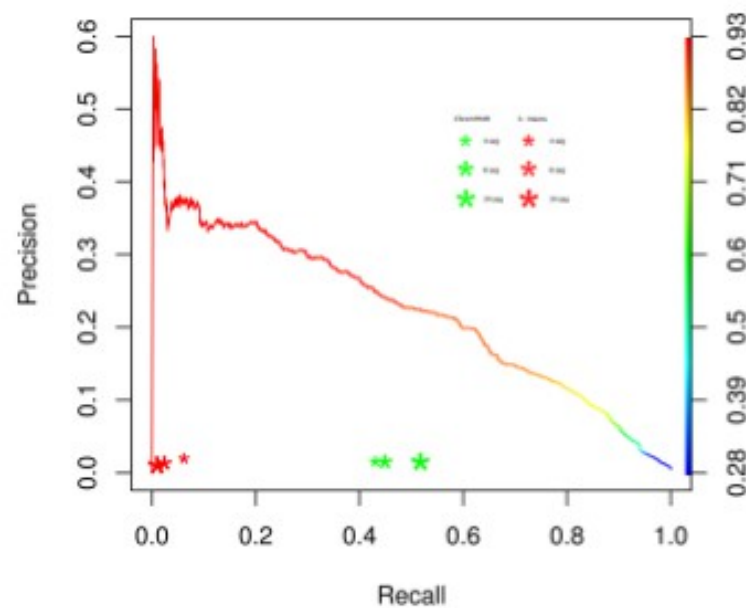
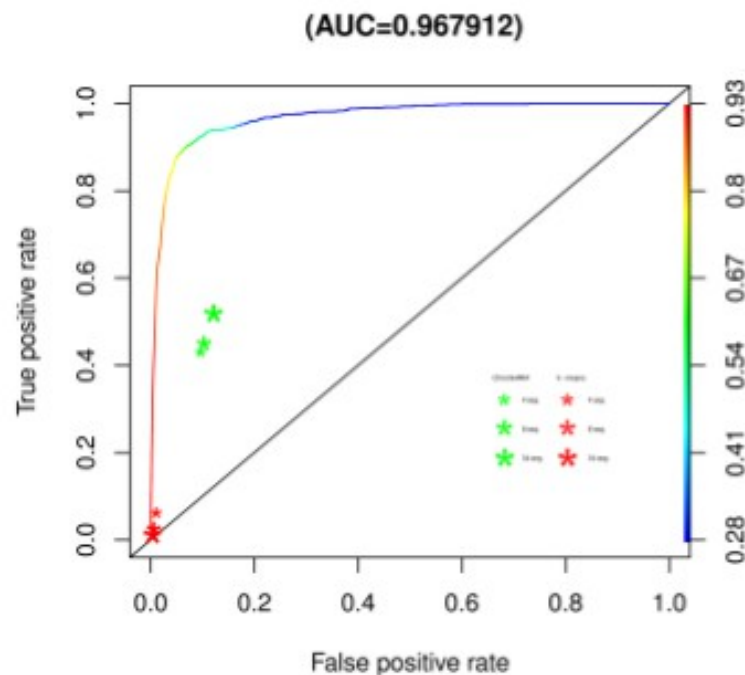


$$P(\text{contact}(X, Y)) = f\left(\sum_{i < k < j} \text{len}(X_k) \cdot \gamma_k\right)$$



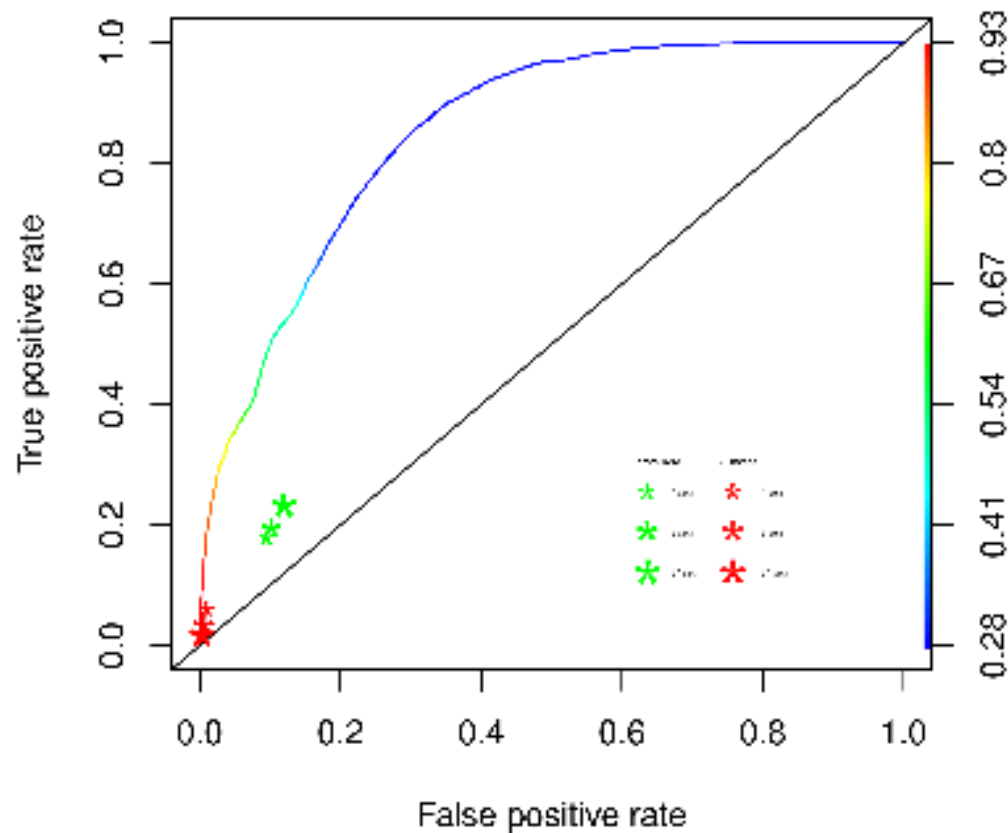
BN classifier can predict boundaries

- Using BN classifiers trained on modENCODE data, we can predict position of boundary elements
- This method outperforms HMMs and clustering of histone modification data

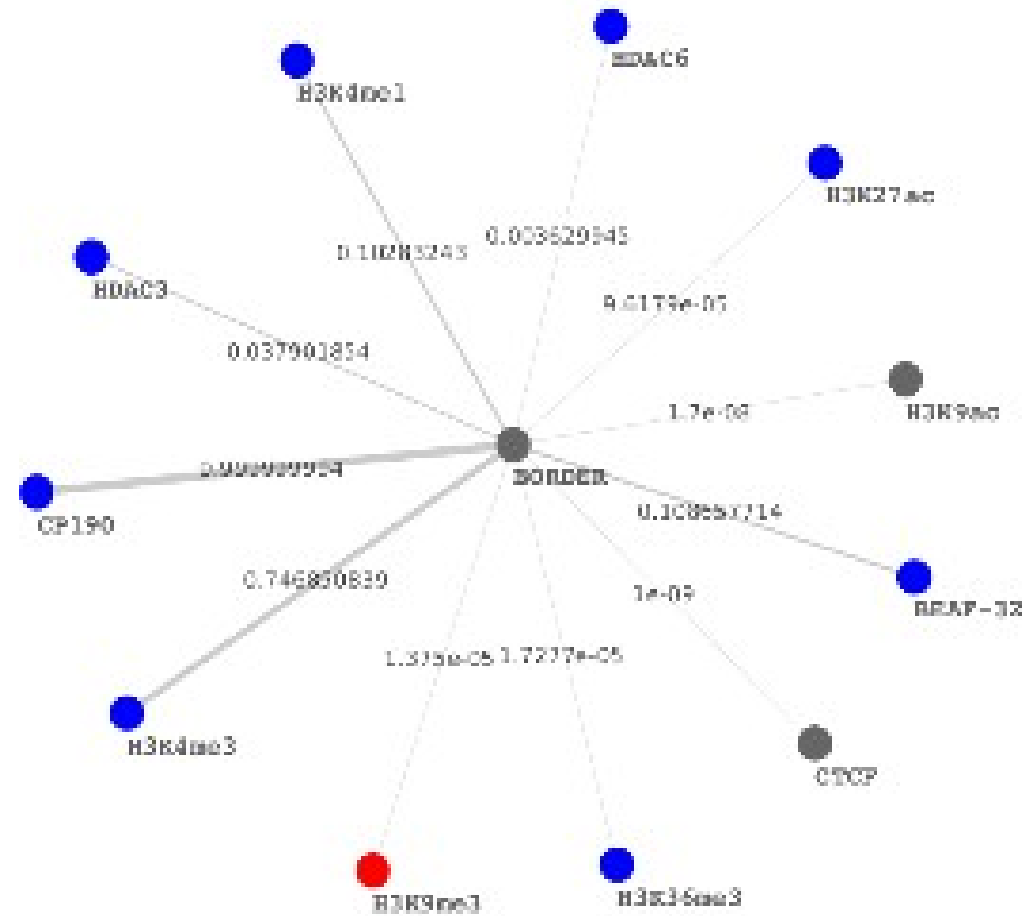


Predictions make sense, and the model brings new information

(AUC=0.853232)



DamID validation results



Impact of signals

Podsumowanie

- Używając zarówno cech sekwencyjnych jak i modyfikacji histonów można dobrze przewidywać lokalizację miejsc funkcjonalnych DNA
- W przypadku cech sekwencyjnych może nam pomóc porównanie spokrewnionych gatunków
- Ta metoda działa nie tylko dla enhancerów ale też dla miejsc izolatorowych