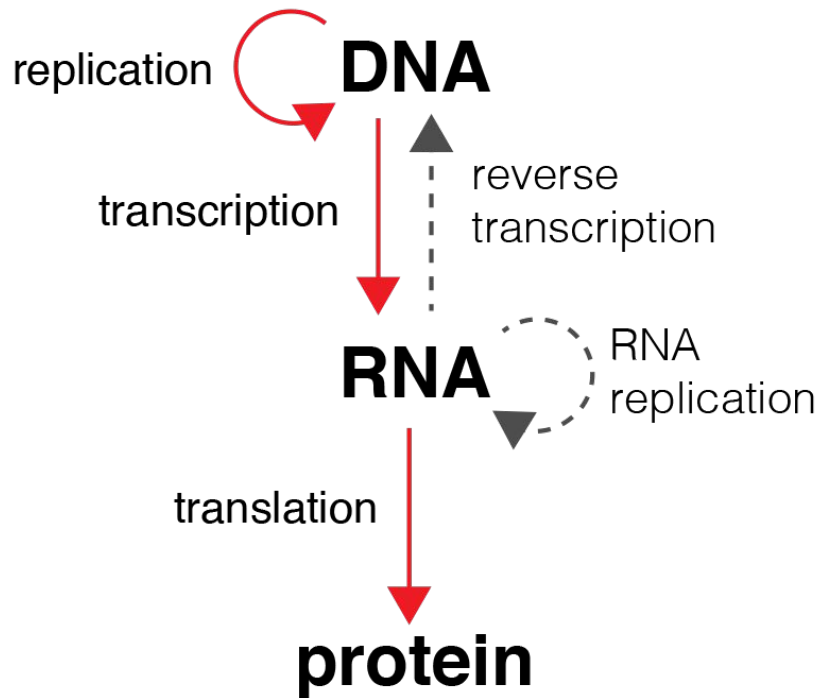


應用資料探勘探討轉錄因子與選擇性剪切的潛在模式

吳旻昇、林子傑

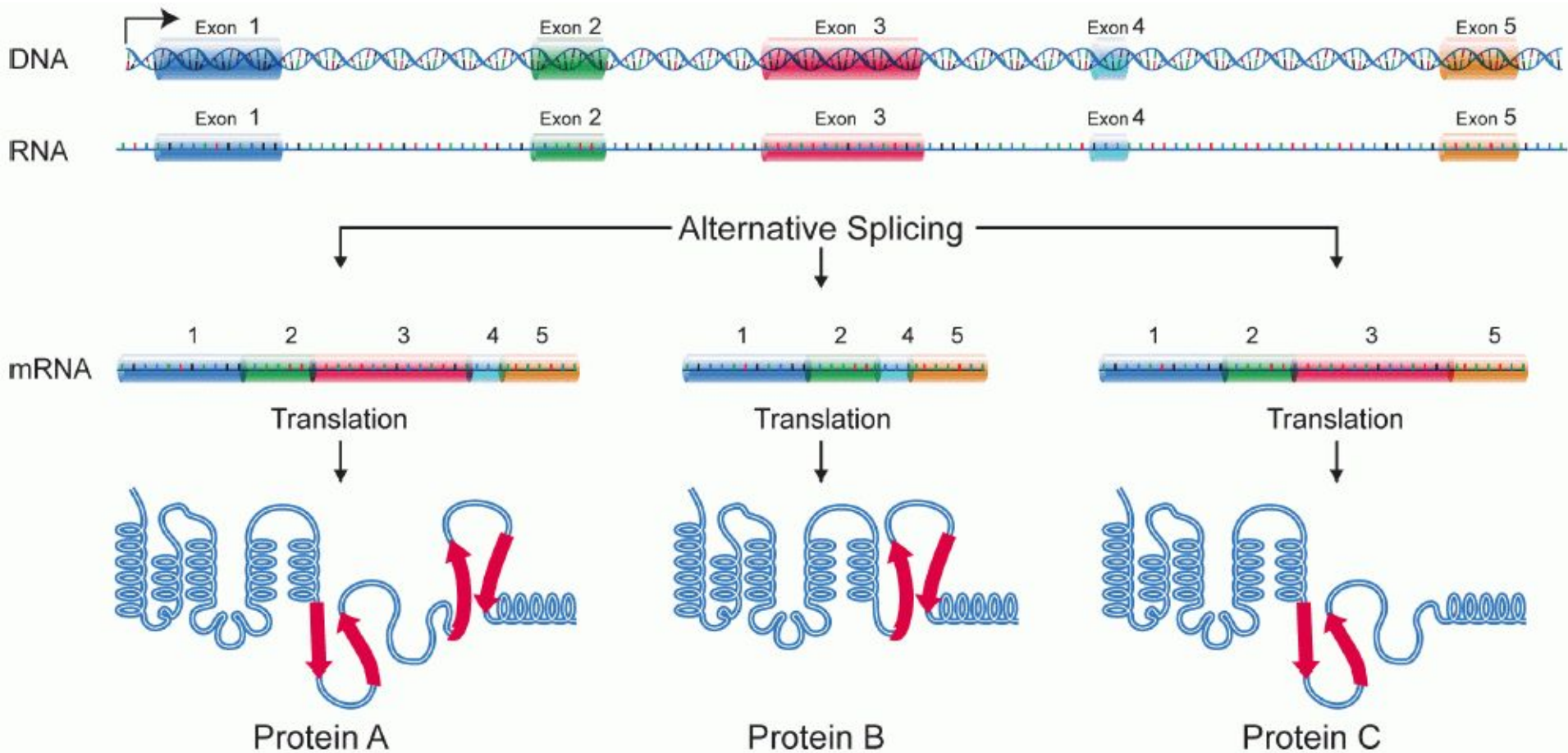
分子生物學的中心教條

- DNA (遺傳訊息) →
RNA (抄本) →
蛋白質 (生理功能)
- 高等生物有非常複雜的生理功能
- 基因數量與複雜程度沒有絕對的正相關

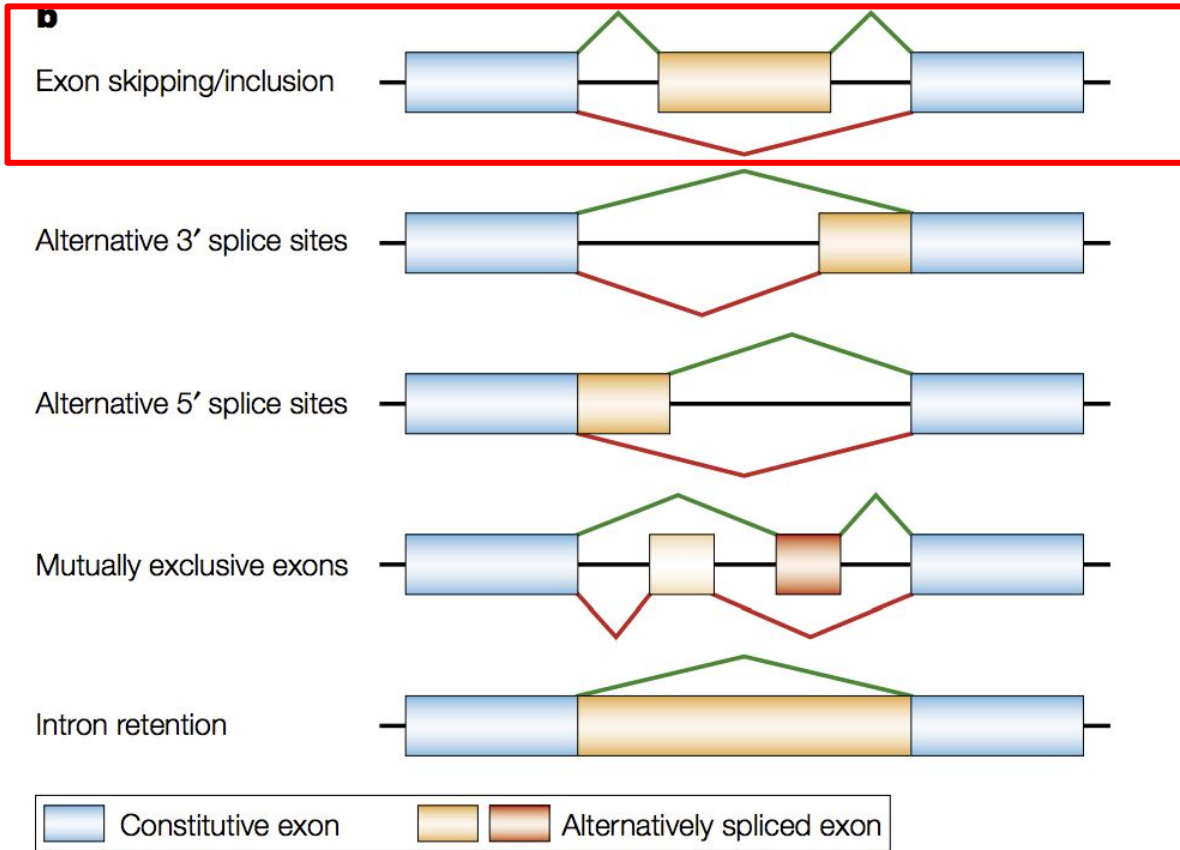


選擇性剪切 (Alternative Splicing)

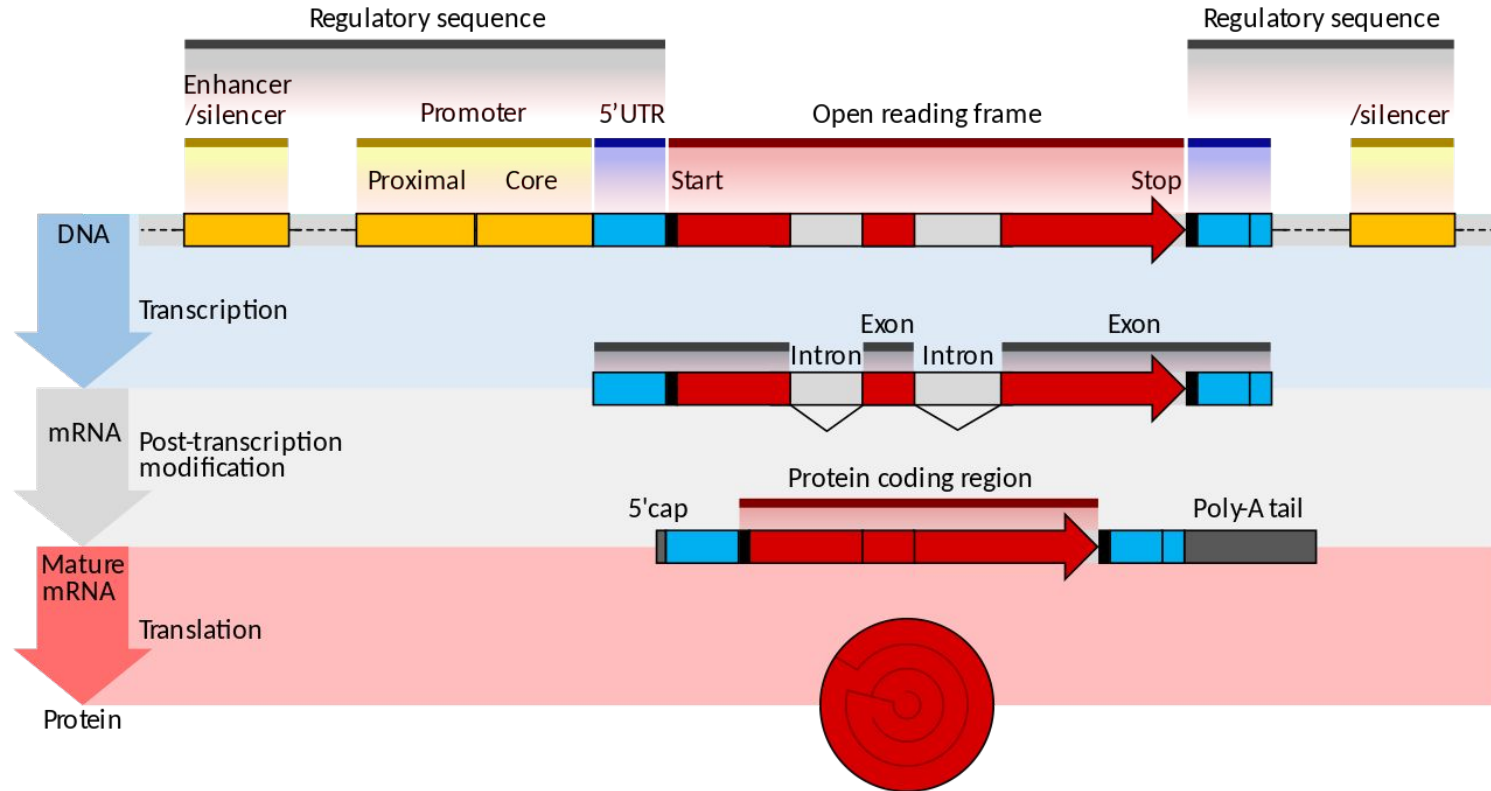
- 在真核生物中並非基因的全部都會轉譯 (translate) 成蛋白質
- Exon (外顯子)
 - 最後會轉譯成蛋白的元件
- Intron (內含子)
 - 在轉錄 (transcription) 後的修飾過程中被剪接去除 (mRNA splicing)
- 透過剪切不同的Exon與Intron組合就可以製造出不同的蛋白



Modes of Alternative Splicing

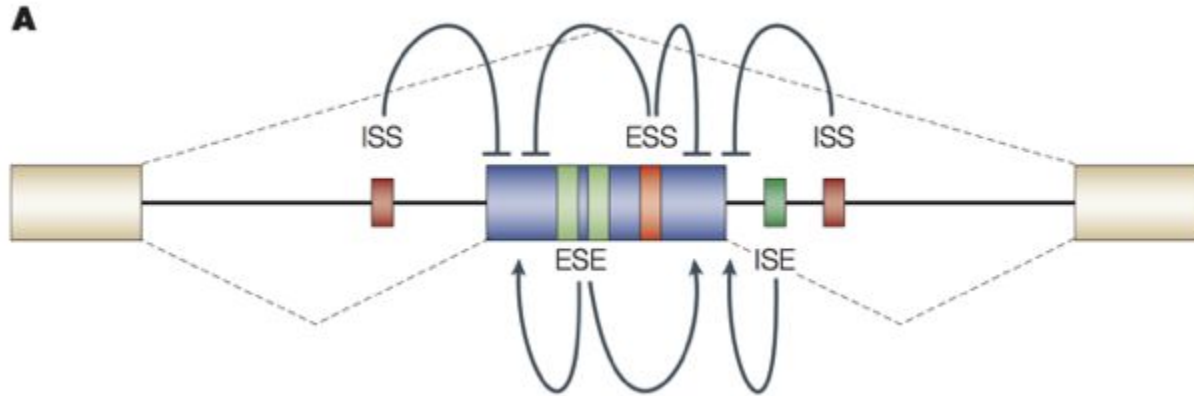


Gene structure



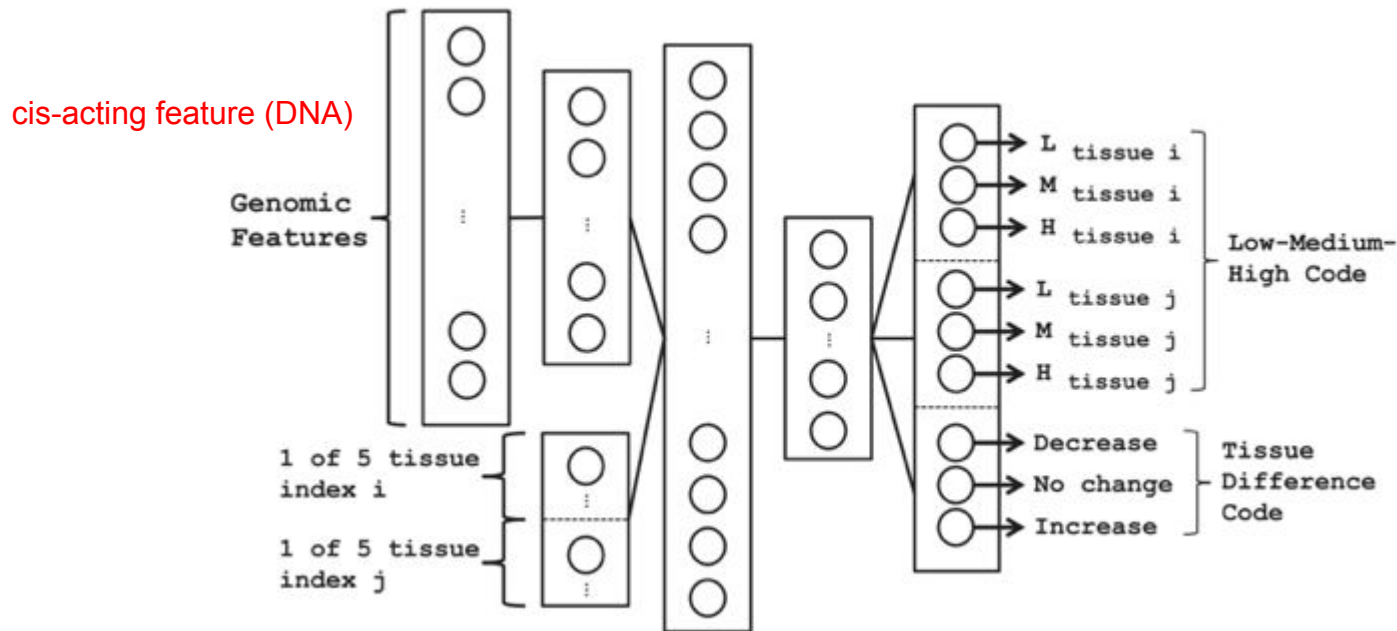
Regulation of Alternative Splicing

- *cis*-acting elements
 - Exon splicing enhancer (ESE)
 - Exon splicing silencer (ESS)
 - Intron splicing enhancer (ISE)
 - Intron splicing silencer (ISS)
- *trans*-acting elements
 - Alternative Splicing Factor
 - Transcription Factor



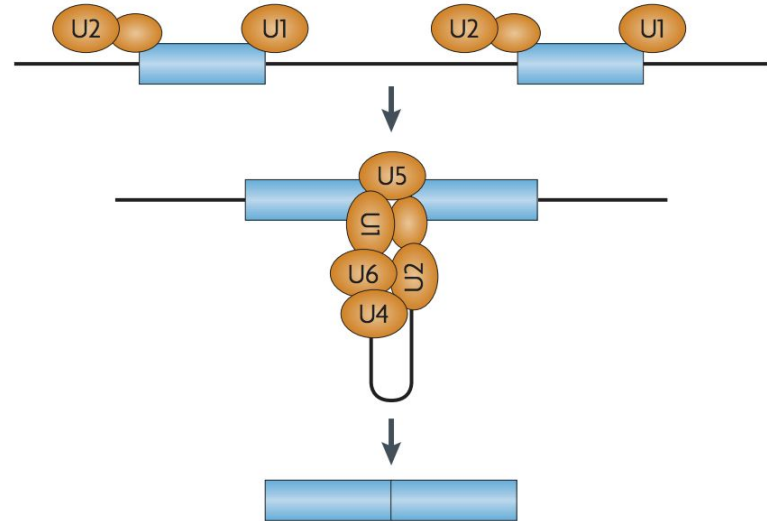
Matlin, A. J., Clark, F., & Smith, C. W. (2005). *Nature reviews. Molecular cell biology*, 6(5), 386.

深度學習預測選擇性剪切



Regulation of Alternative Splicing

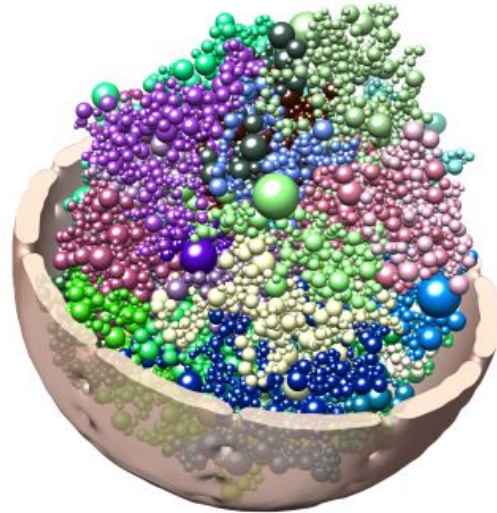
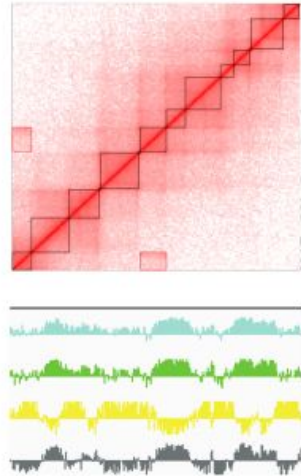
- *cis*-acting elements
 - Exon splicing enhancer (ESE)
 - Exon splicing silencer (ESS)
 - Intron splicing enhancer (ISE)
 - Intron splicing silencer (ISS)
- *trans*-acting elements
 - Alternative Splicing Factor
 - Transcription Factor

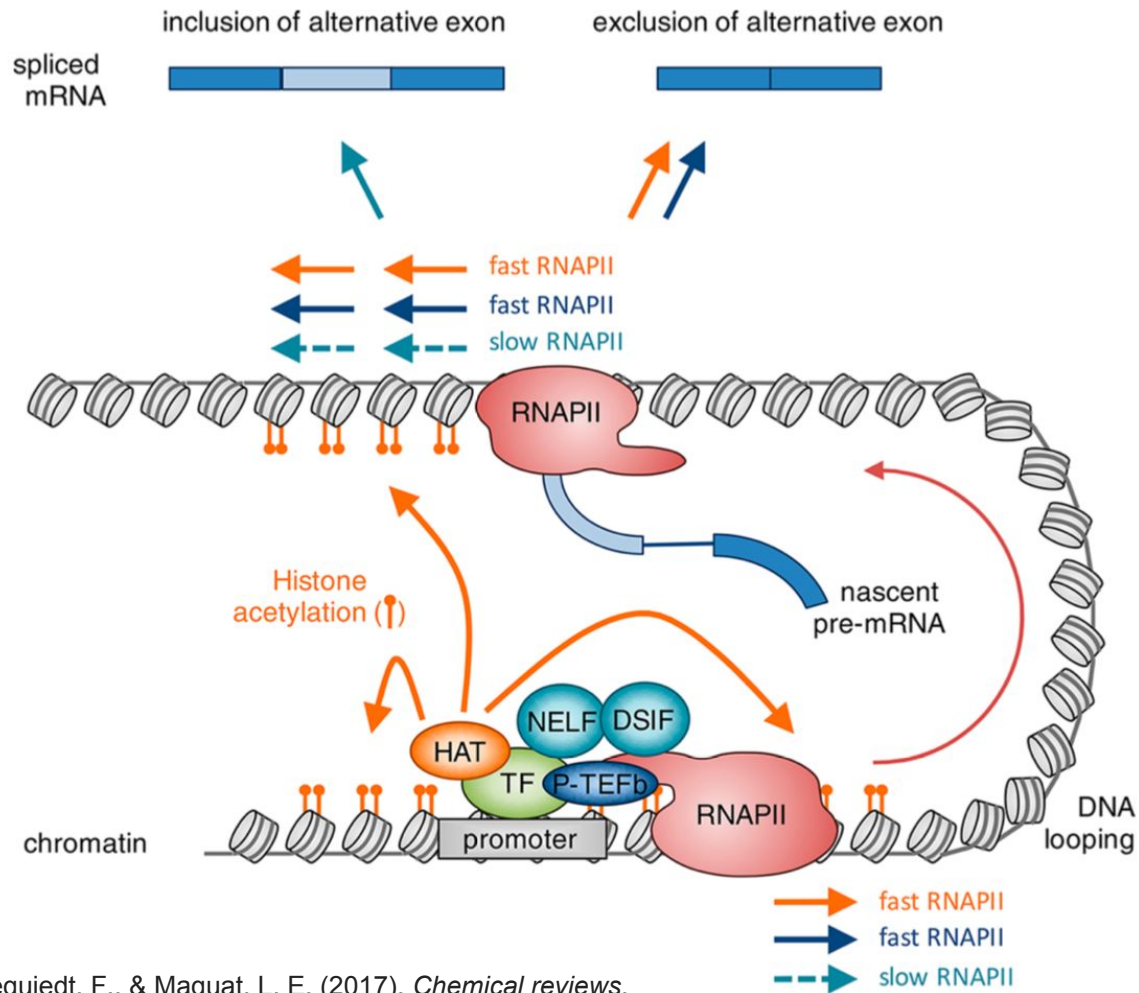


Keren, H., Lev-Maor, G., & Ast, G. (2010).
Nature Reviews Genetics, 11(5), 345.

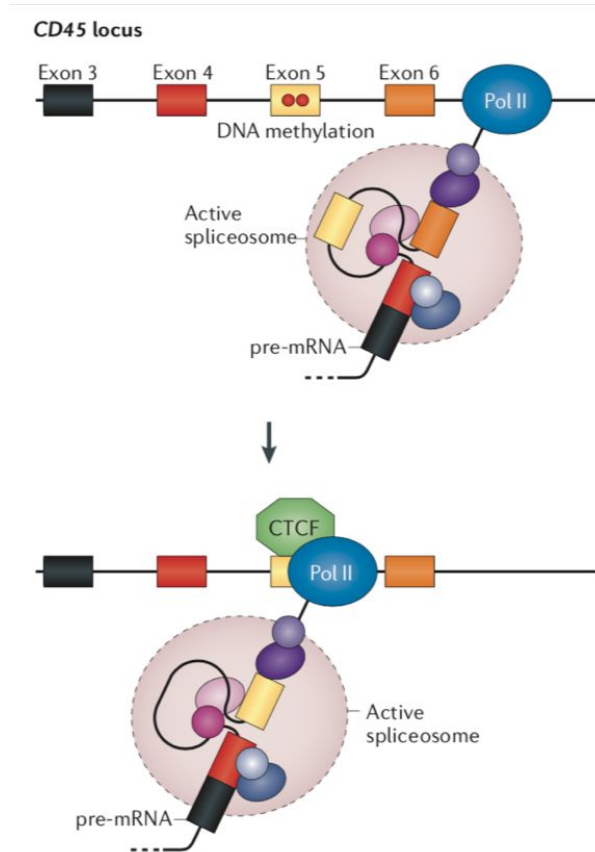
Distal DNA Element can influence AS

- Chromatin is organized in a 3D struction
- Distal DNA element may have tightly connection
- Hi-C Sequencing reveal complex 3D chromatin structure in nucleolus





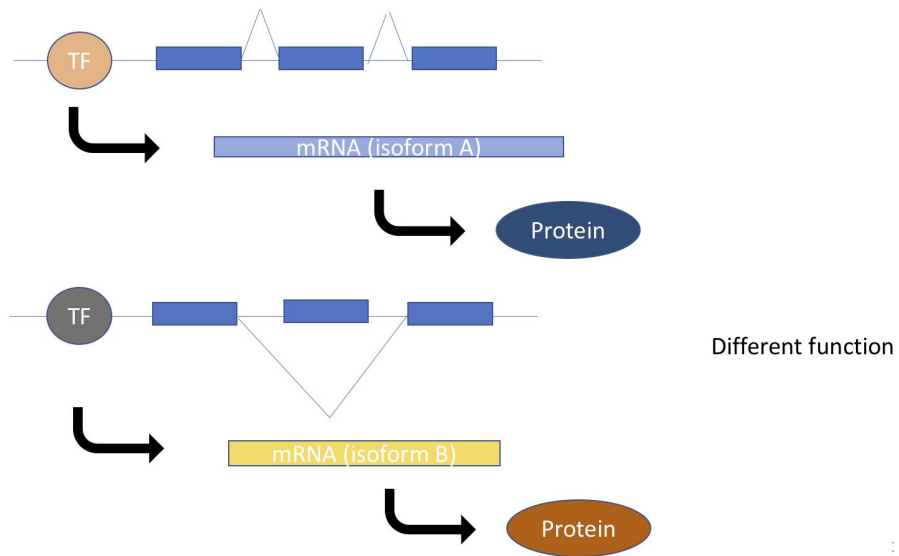
CTCF Promote Alternative mRNA



Ong, C. T., & Corces, V. G. (2014). *Nature reviews. Genetics*, 15(4), 234.

目標

- Gene expression is regulated by TF. How about alternative splicing?



資料來源

- Drosophila melanogaster (果蠅)
- Genomic data
 - Ensembl BDGP6
 - https://www.ensembl.org/Drosophila_melanogaster/Info/Index
- Annotation
 - BDGP6.90
- FAIRE-Seq & RNA-Seq from NCBI GEO Database
 - GES38727, GSE40739 ,GSE62558
 - <https://www.ncbi.nlm.nih.gov/geo/>

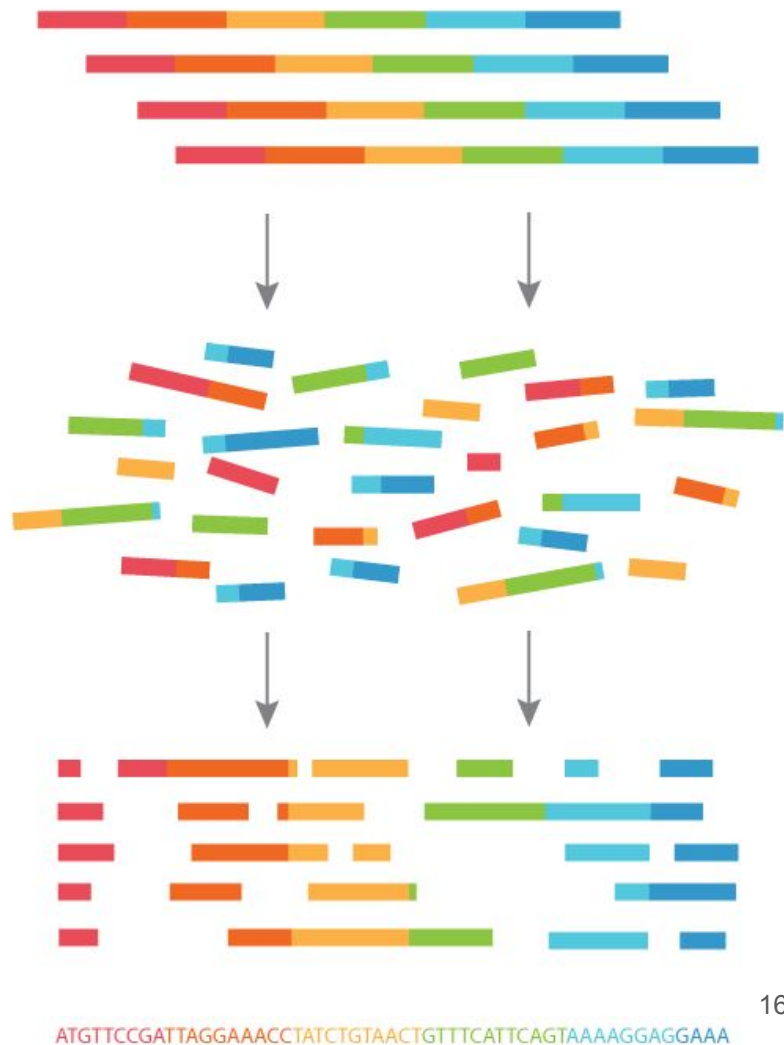


Next Generation Sequencing

- Procs
 - 高通量 (High-Through)
 - 單位成本低廉
- Concs
 - 短讀長 (short reads)
 - 相對低的準確率
- 需要有效率演算法來處理數據

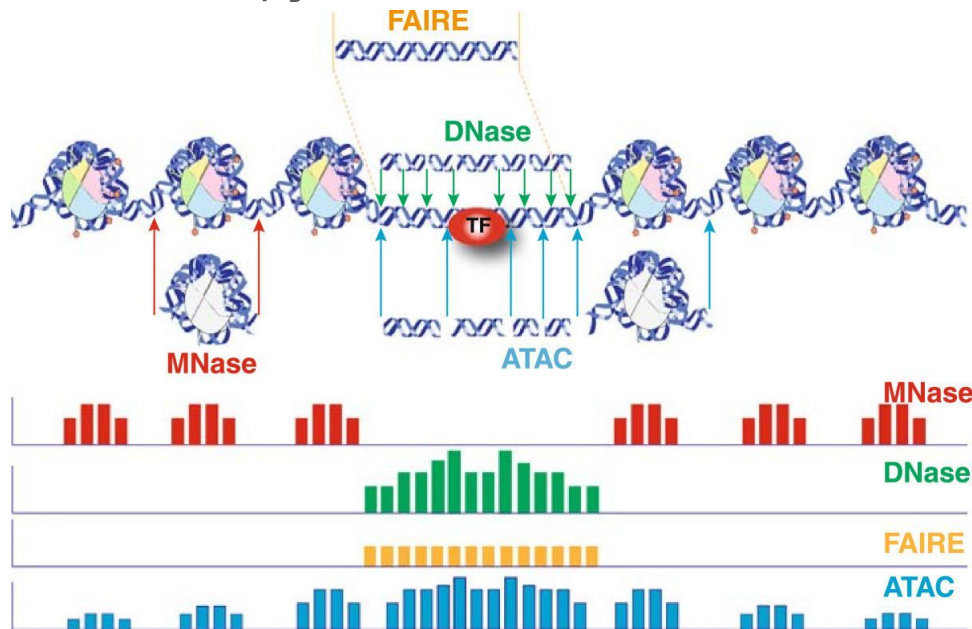
illumina®

Roche
454
SEQUENCING



FAIRE-Seq

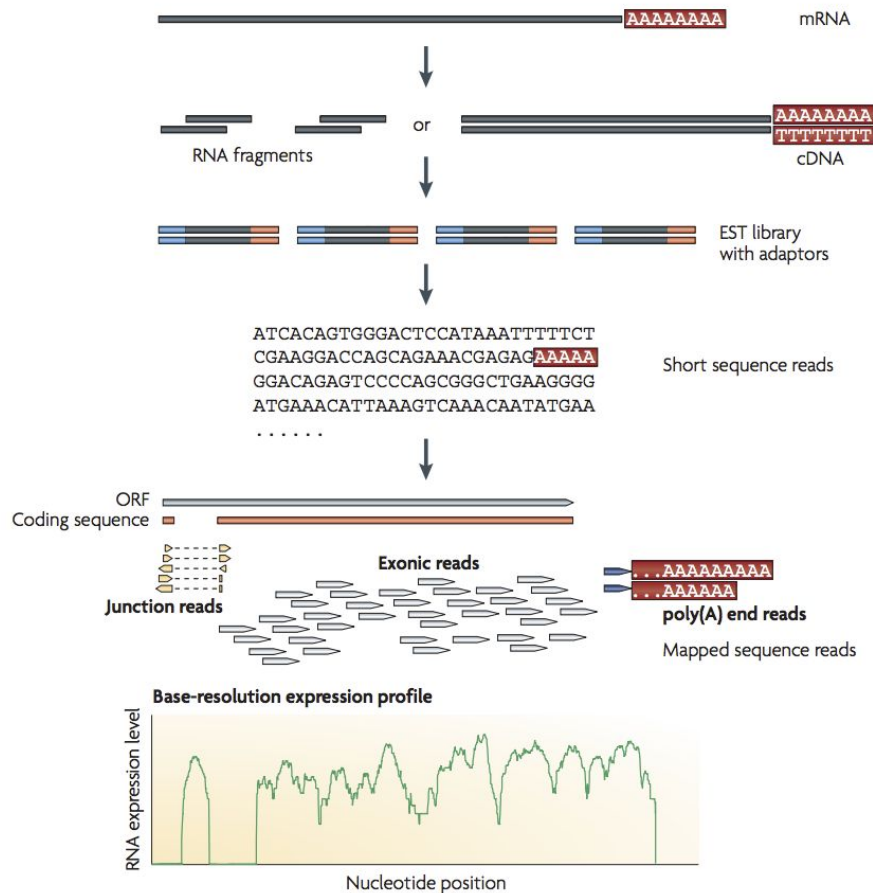
- 與轉錄因子結合的DNA區域，染色質會呈現打開的狀態 (chromatin-accessible)
- 利用formaldehyde固定具有Protein-DNA interaction的區域
- 分離chromatin-accessible的DNA



Tsompana, M., & Buck, M. J.
(2014). *Epigenetics & chromatin*,
7(1), 33.

RNA-Seq

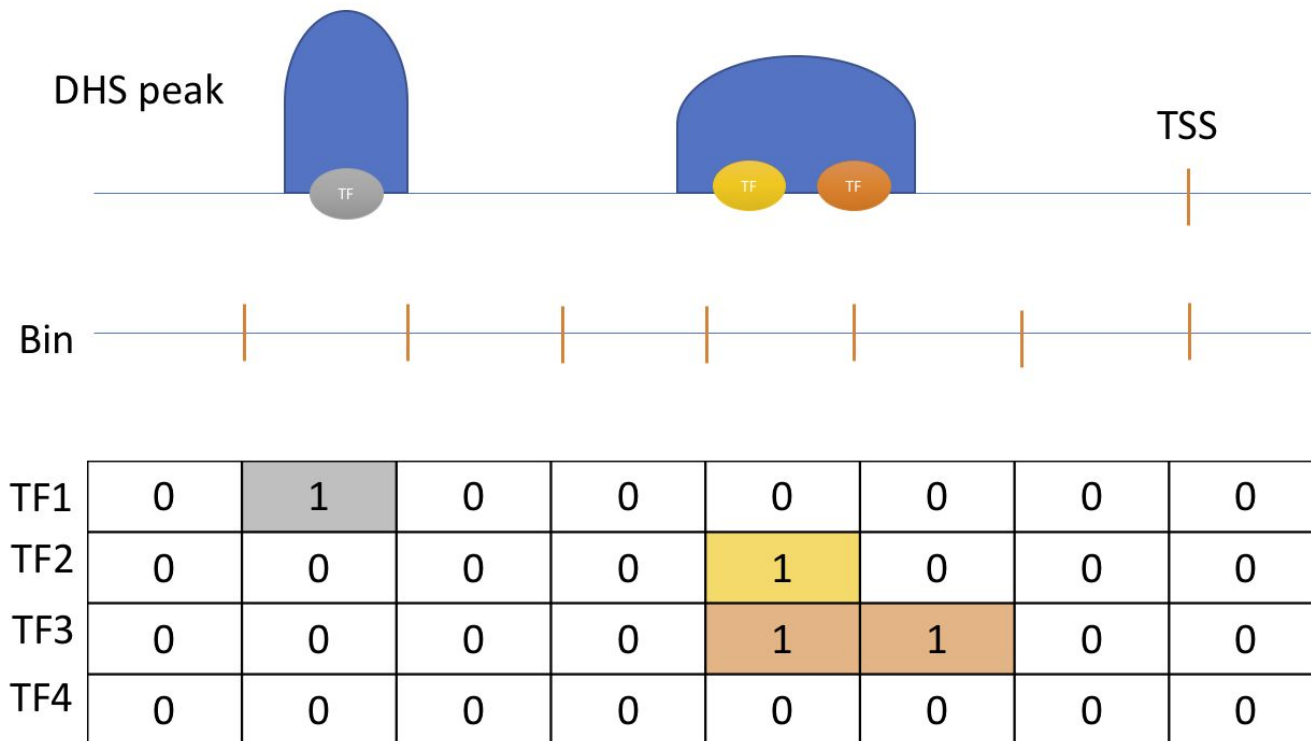
- 利用定序來定量以及定性基因的表現量
- 直接定序不同的 Isoform
- 可用來研究選擇性剪切



資料前處理

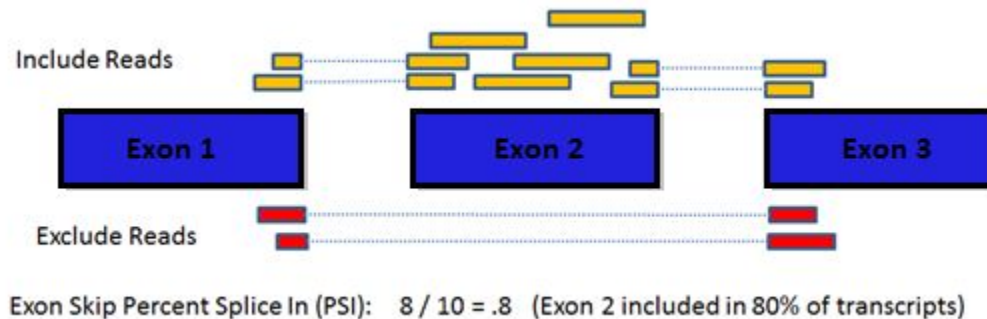
- RAW NGS data
 - Trimmomatic (QC)
- RNA-Seq
 - HISAT2 (Short Reads Mapping)
 - CATANA (AS events from annotation)
 - MISO (Quantify AS events)
- FAIRE-Seq
 - Bowtie2 (Short Reads Mapping)
 - MACS2 (Peak Calling)
 - CIS-BP (Motif Database)
 - RAST: Matrix-SCAN (Find TF Binding sites in Peak)

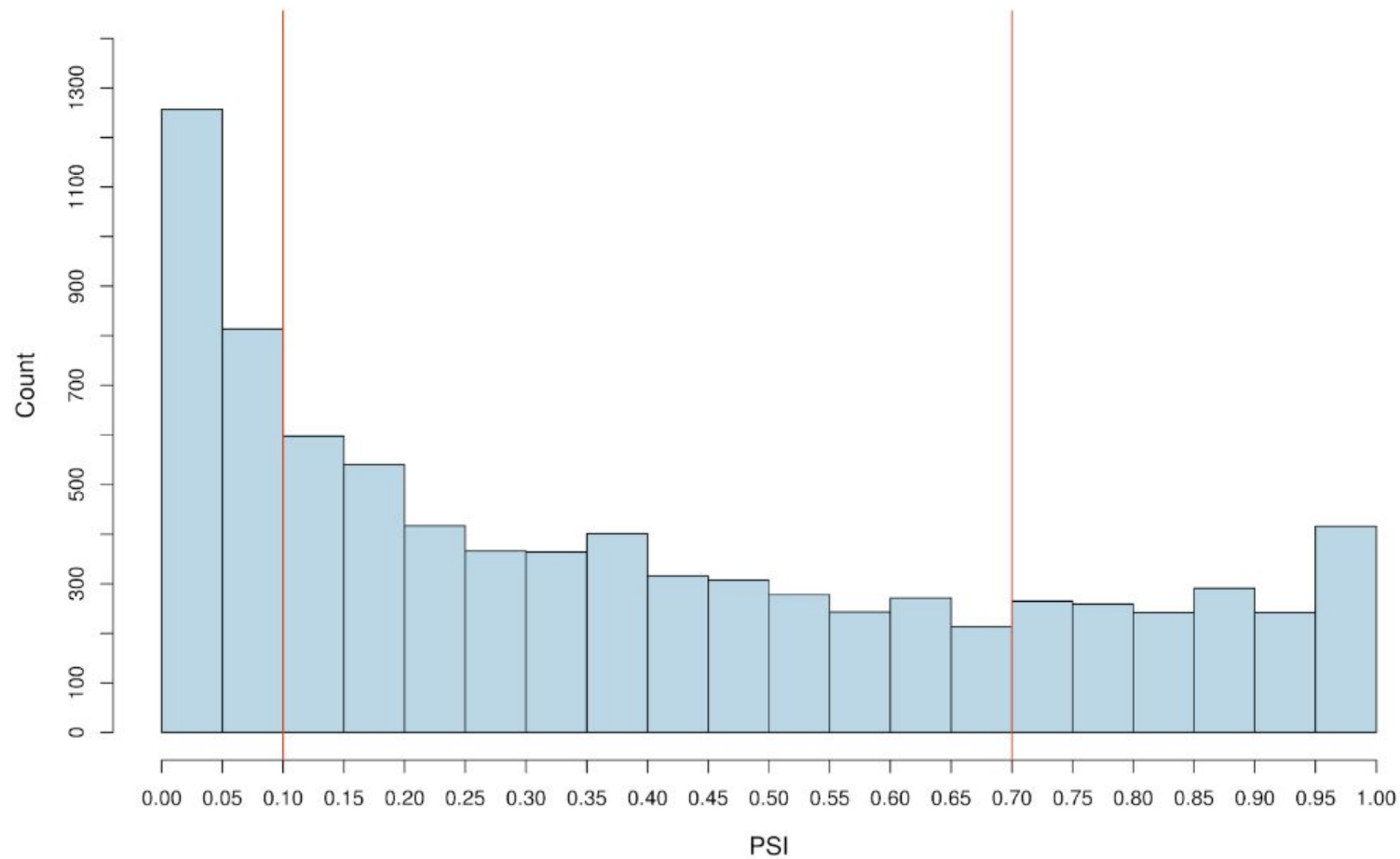
Feature extraction



Targets

- 只取一個基因的第一個Exon Skipping Event
- 利用MISO計算PSI值
- 依PSI值分成H/L兩類
 - $H > 0.7$, $L < 0.1$





Spliced: 2071

Retained: 1715

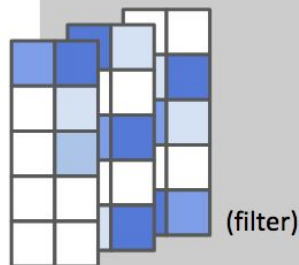
Deep Learning Model

- Activation Function: ELU with Batch Normalization
- CNN
 - Conv1D -> Max Pooling -> Flatten -> FC1 -> FC2 -> Output (Softmax)
 - Dropout (rate): 0.5
- DNN
 - Flatten -> FC1 -> FC2 -> FC3 -> FC4 -> Output (Softmax)
 - Dropout (rate)
 - FC1,2: 0.5
 - FC3,4: 0.2

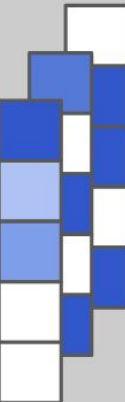
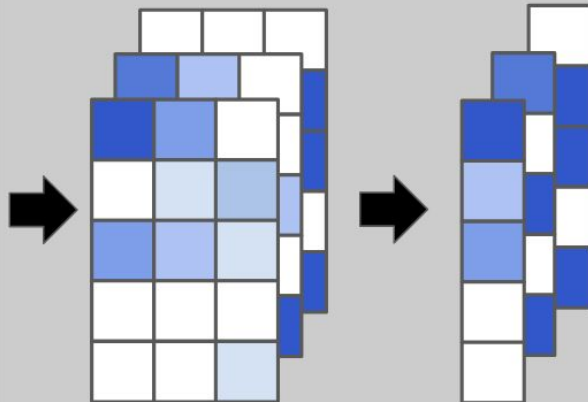
→

T	T	F	F
F	F	F	T
F	T	T	F
F	F	F	F
F	F	T	T

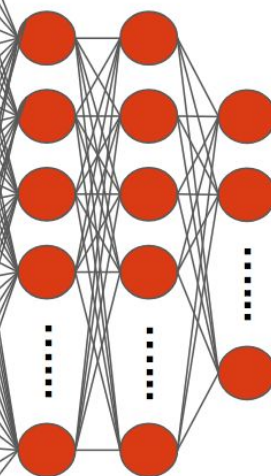
Matrix



(CNN)



(DNN)



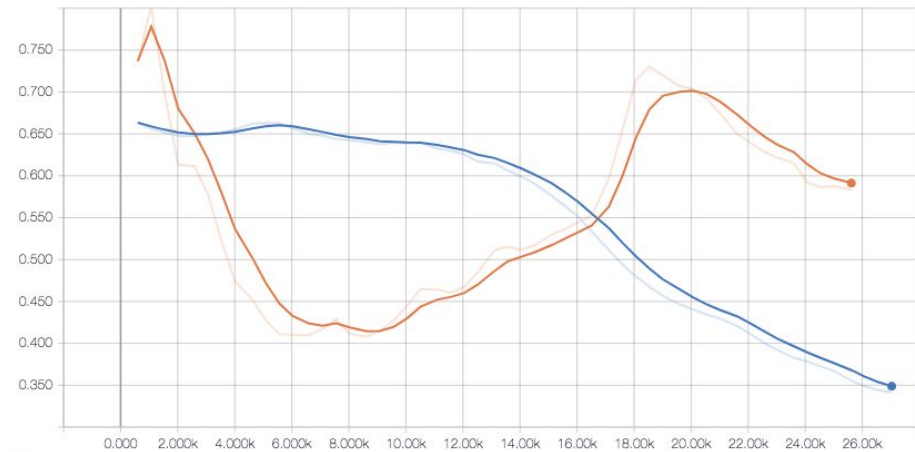
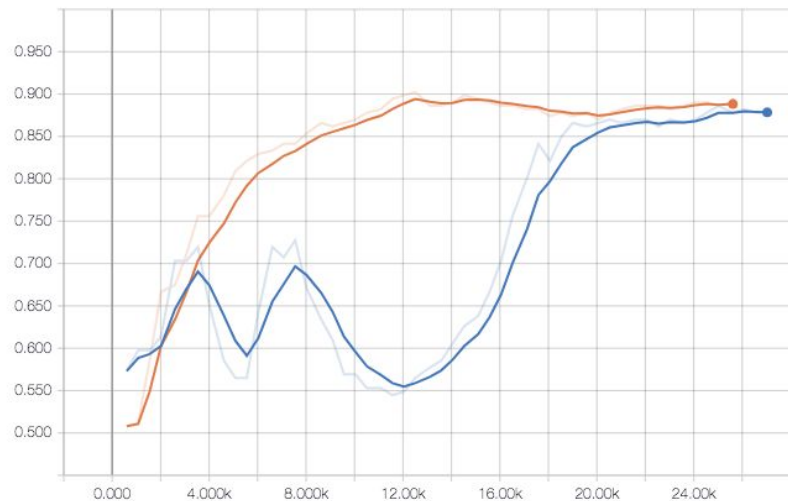
Splice
or
Non-splice

Deep Learning Model

- Mini Batch training
 - Batch size: 128
- Initial learning rate: 0.001
- L2 regularization
 - Wight decay
 - CNN: 0.01
 - DNN: 0.001
- Optimizer: Adam
- Train/Validation/Test: 2832/246/708

Training

- Orange: DNN, Blue: CNN

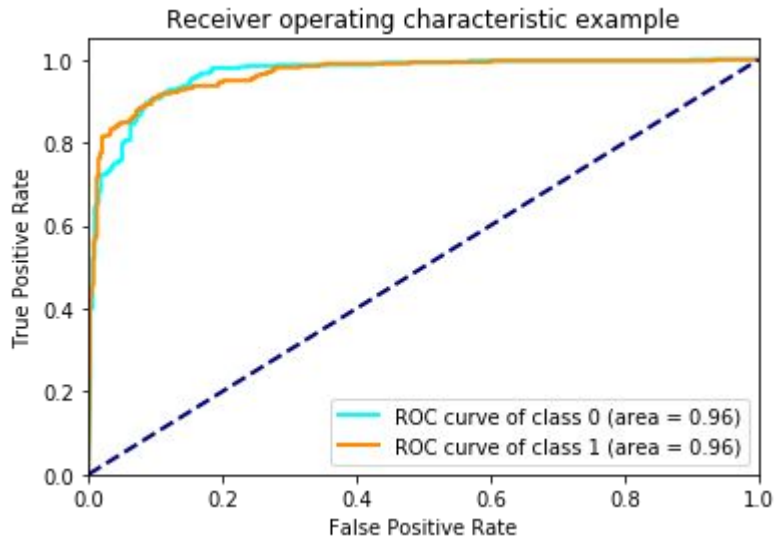
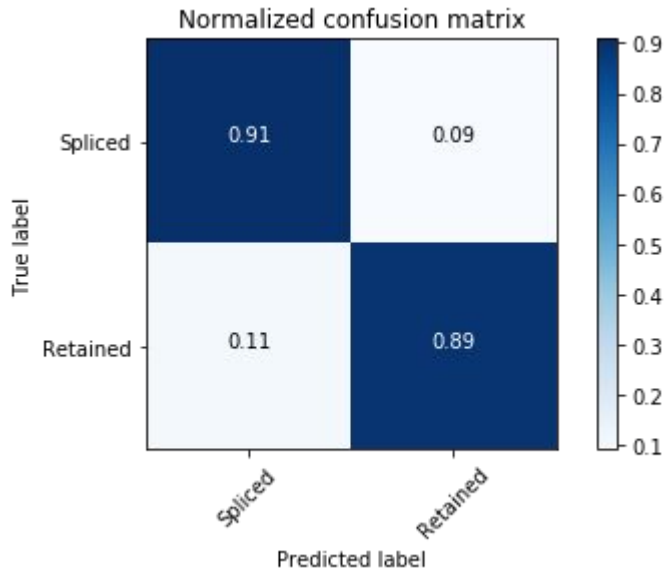


10-fold Cross Validation

- CNN
 - Loss: 0.310 +- 0.032
 - Accuracy: 87.13 +- 1.59%
- DNN
 - Loss: 0.631 +- 0.104
 - Accuracy: 90.16 +- 0.96%

模型評估

- 從PSI轉換成H/L類別 (Spliced/Retained) 時沒有平衡兩個類別
- 所以分類的混淆矩陣兩個類別稍有不平衡



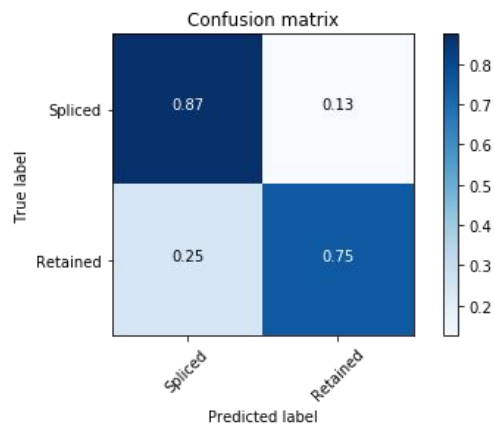
XGBoost

- Deep Learning很強大, 但是不能夠解釋究竟那個Feature (TF)是影響選擇性剪切的關鍵
- XGBoost
 - A Scalable Tree Boosting System
 - 實際上除了Boosting, XGBoost同時也做了Bagging
 - 同時利用一階導數及二階導數來優化權重
 - 可以利用Information Gain來找出Feature Importance

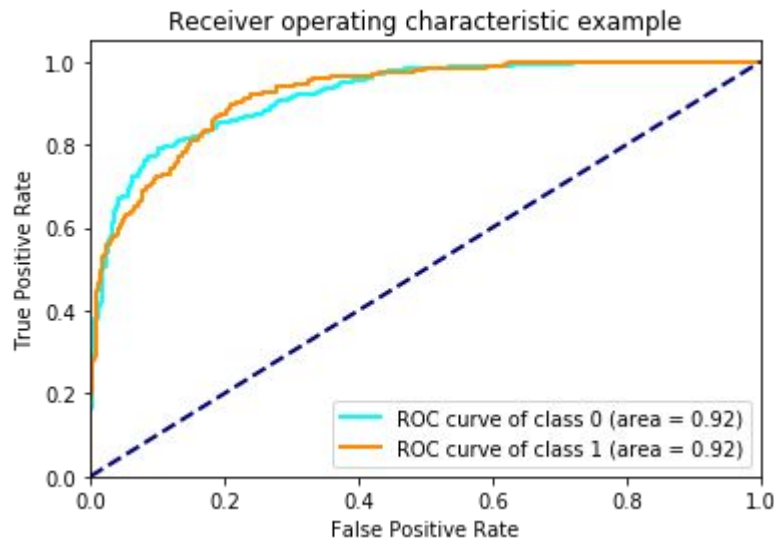
XGBoost

- N_estimators: 500
- Learning_rate: 0.1
- Max_depth: $\sqrt{n_{TF}} = \sqrt{307}$
- Train/Test split: 0.8/0.2
- Performance
 - Accuracy: 81.66%
 - AUROC: 0.92

模型評估



稍微Unbalance

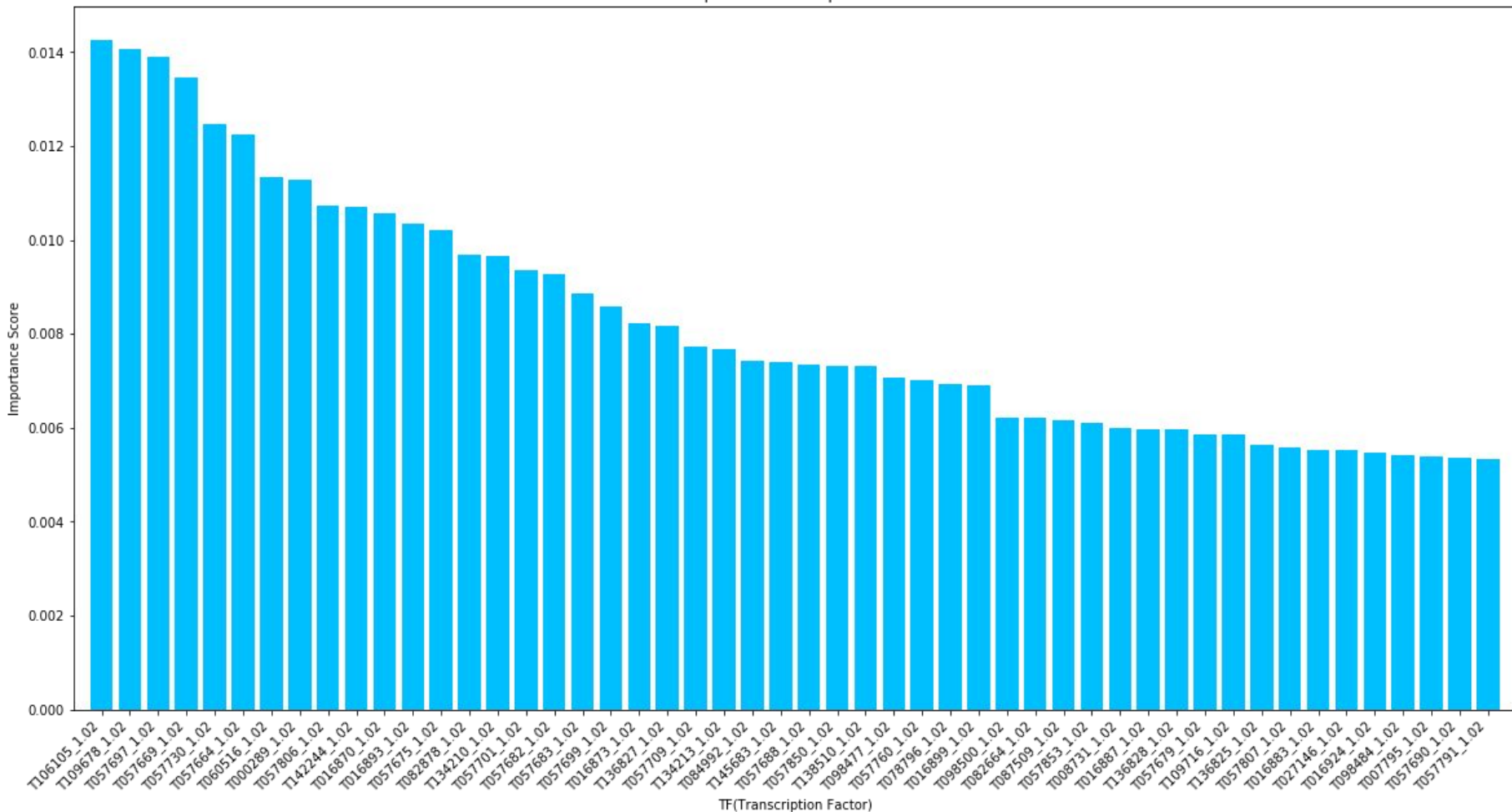


5-fold Cross-Validation

- Accuracy: 83.44 +- 0.93%
- AUC: 0.924 +- 0.009

- 因爲TF在每個Fold中的Feature Importance有時變化很大
 - 利用5個Fold的Rank做平均
 - 避免訓練時 Stochastic的影響

Top 50 Feature importances



Top-10 TF

1. T106105_1.02 (0.014262) (Hsf)
2. T109678_1.02 (0.014078) (Adf1)
3. T057697_1.02 (0.013909) (lola)
4. T057669_1.02 (0.013463) (hb)
5. T057730_1.02 (0.012457) (jim)
6. T057664_1.02 (0.012262) (Cf2)
7. T060516_1.02 (0.011338) (rn)
8. T000289_1.02 (0.011273) (br)
9. T057806_1.02 (0.010726) (CTCF, CCCTC-binding factor)
10. T142244_1.02 (0.010704) (Mad)

Previous study

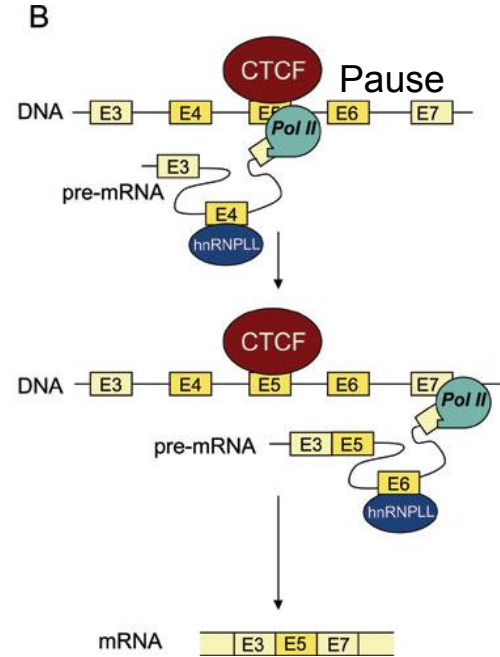
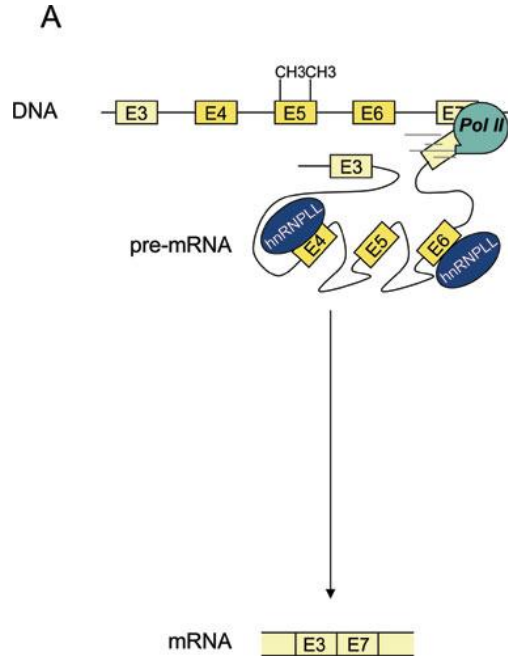
ARTICLE

doi:10.1038/nature10442

CTCF-promoted RNA polymerase II pausing links DNA methylation to splicing

Sanjeev Shukla¹, Ersen Kavak^{2,3}, Melissa Gregory¹, Masahiko Imashimizu⁴, Bojan Shutinoski¹, Mikhail Kashlev⁴, Philipp Oberdoerffer¹, Rickard Sandberg^{2,3} & Shalini Oberdoerffer¹

Mammalian-model (哺乳類)



Shukla S et al. (2011) Nature. 479:74–79

Kornblihtt AR (2012) Cell Res. 22(3):450-2

結論

- 透過深度學習以及XGBoost發現轉錄因子與選擇性剪切存在潛在的Pattern
- 轉錄因子極有可能是調控選擇性剪切的因子之一
 - CTCF已經被證實能夠調控Exon Skipping (哺乳類)
 - CTCF在果蠅都發現可能存在調控 AS功能
 - TF調控AS的機制具有廣泛性, 因此可能在很早期的真核生物就已經存在 (待更多證據)
- 沒有少數重要的TF, 所以Combinatorial Effect是需要被考慮的
- 資料探勘協助發現潛在的生物問題與機制

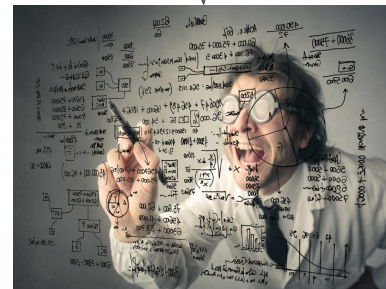


生物問題



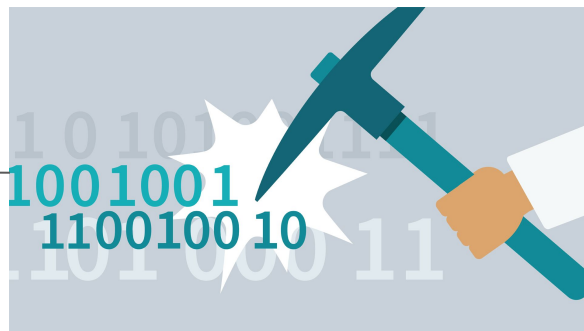
生物實驗

產生資料與結果



資料科學家

重新分析



資料探勘

Data-Driven Biological Research

生物學家

給予生物學家 Insight