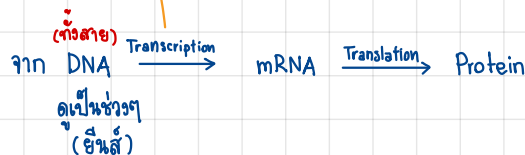


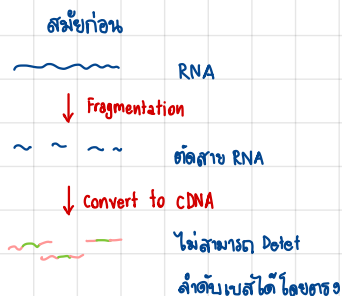
# Xpore

## Problem statement

### ▷ Central Dogma



### ▷ RNA Sequencing



Nanopore ขนาดเล็กสุดคือ MinION

ปัจจุบัน (ใช้ Nanopore)

สามารถทำ direct RNA sequencing ได้

เช็ค RNA modification

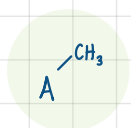
หลักการทํางาน



โมเลกุลแต่ละตัวจะทำให้เกิดความต้านทานแตกต่างกัน และใช้ ML ในการอ่านสัญญาณไฟฟ้า เป็นลำดับเบส (Base calling)

### ▷ RNA modifications

เช่น m<sup>6</sup>A



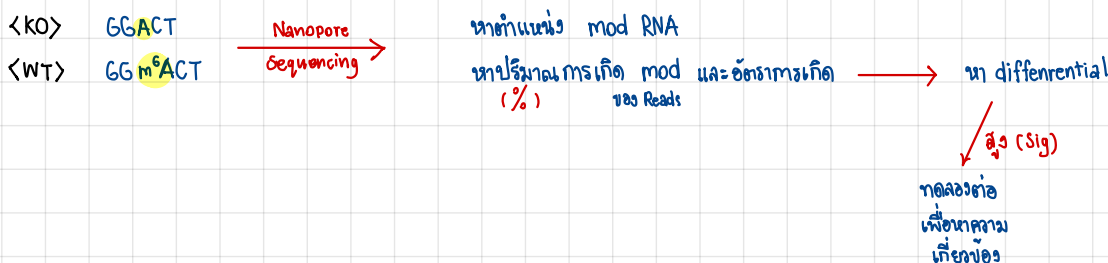
Sample

KO - knock out : เป็น sample ที่ knock out ยีนส์ที่ผลิต m<sup>6</sup>A ออก

WT - wild type : ปกติ

เปรียบเทียบ KO-WT เพื่อหา modification rate

### ▷ Research Objectives



### - Signal-level modification detection methods

EpiNano / MINES (software)

- สามารถ detect ตำแหน่ง m<sup>6</sup>A ได้
- จำเป็นต้องมี data ใน training (Supervised)

Tombo

- สามารถ detect ตำแหน่งของ mod ได้จากประมาณ
- ไม่ต้องใช้ data ในการ train (Unsupervised)

Xpore อยู่ประเภทนี้

## ① Data collection and Preparation (จาก Xpore)

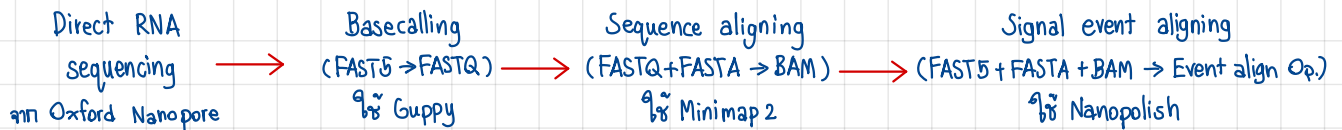
FAST5 file คือไฟล์ สัญญาณไฟฟ้า (pA) ถูกเก็บไว้ในรูปแบบ HDF5 (binary)

FASTQ file คือไฟล์ ลำดับเบสในรูปแบบ text

FASTA file คือไฟล์ ลำดับเบส cdna (rna → dna)

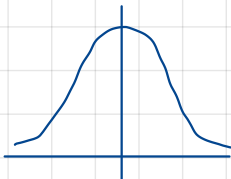
BAM/SAM คือไฟล์ ที่จัดเรียง FASTQ และ FASTA  
Binary text

Nanopore preprocessing



## ② Bayesian [Multi-Sample] Gaussian Mixture Modelling

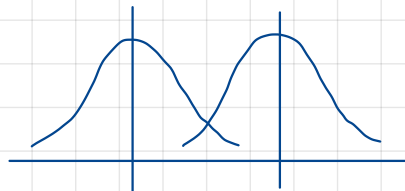
Gaussian (normal distribution)



$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$$

⇒ รู้ assumptions ว่า data ถูกสร้างอย่างไร จึงสามารถสร้าง data ขึ้นได้ (Generative AI ตัวยกรก)

2 Gaussian Mixture



มี 2 distribution

Frequentis }  
Bayasian } 1 parameter

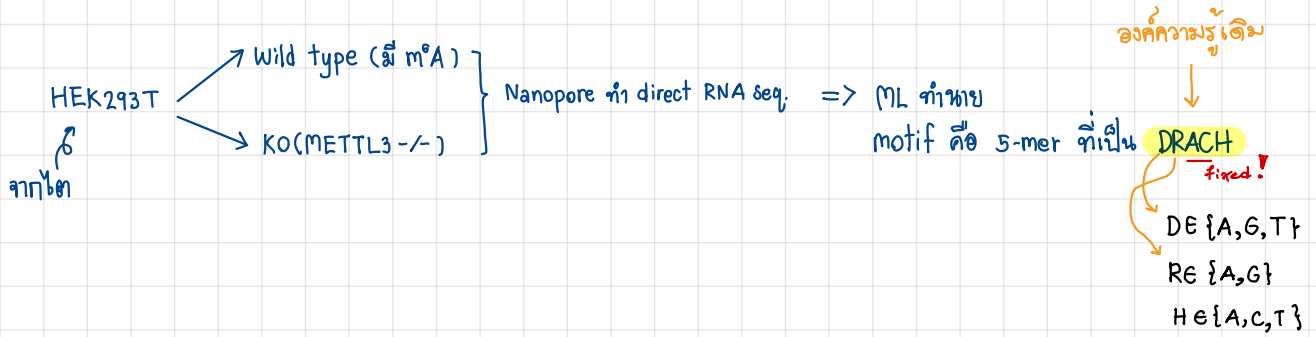
Multi-Sample - เช่น  $ma_1, ma_2, ma_3$  as train

ในงานวิจัยนี้ใช้ GMM เพราะ - มี 2 distribution ที่น่าจะเป็น คือ 1) modified ( $m^*A$ )  
2) unmodified ( $A$ )  
- รู้ prior จาก nanopolish

### ③ Evaluation

Cell line - เซลล์ที่เลี้ยงให้เติบโต และเพิ่มจำนวนได้อ่อนแอ

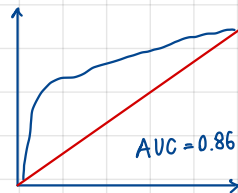
#### ▷ Experiment Setup



#### ▷ Validation: m6A calling ให้ ROC Curve

##### In Bioinformatics

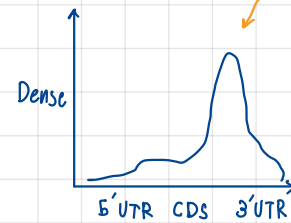
- false positive อาจจะไม่ผิดเสมอไป
- จำเป็นต้องทำ Analysis อื่นๆ เพื่อจะได้ Insight มากขึ้น



- มีความแม่นยำ 86%



- m<sup>6</sup>ACE-Seq + DRACH  
 มี Acc > 95%



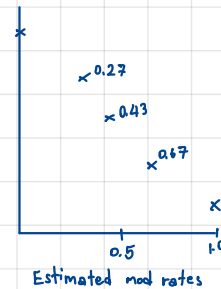
Confirm องค์ความรู้เดิม!

จะเกิด m<sup>6</sup>A และมักจะเกิดช่วงท้ายๆ ของ RNA

#### m6A stoichiometry quantification

RNA Mixture  
 HEK293-KO + HEK293-WT  
 0%  
 25%  
 50%  
 75%  
 100%  
 (m<sup>6</sup>A %)

Modification rates: GGACT



∴ สามารถหา Modification rates ได้

สรุป ML Metrics : ROC curve, Precision-Recall Curve, Accuracy  
 Analysis : Domain-specific evaluation, Effects of the data size

#### ▷ Applicability : Full Dataset

ทำกระบวนการที่ทำให้คนทั่วไปใช้งาน ใช้ได้ง่ายขึ้น  
 เช่น เปรียบเทียบหลายๆแบบ (โต, น้อย, ปอด, ฯลฯ)  
 หรือเปรียบเทียบเป็นคนๆ

#### Keys Takeaway

- Validation - ใช้ metrics ที่เหมาะสม
- Comparison กับของคนอื่น
- Applicability - นำไปใช้ได้จริง