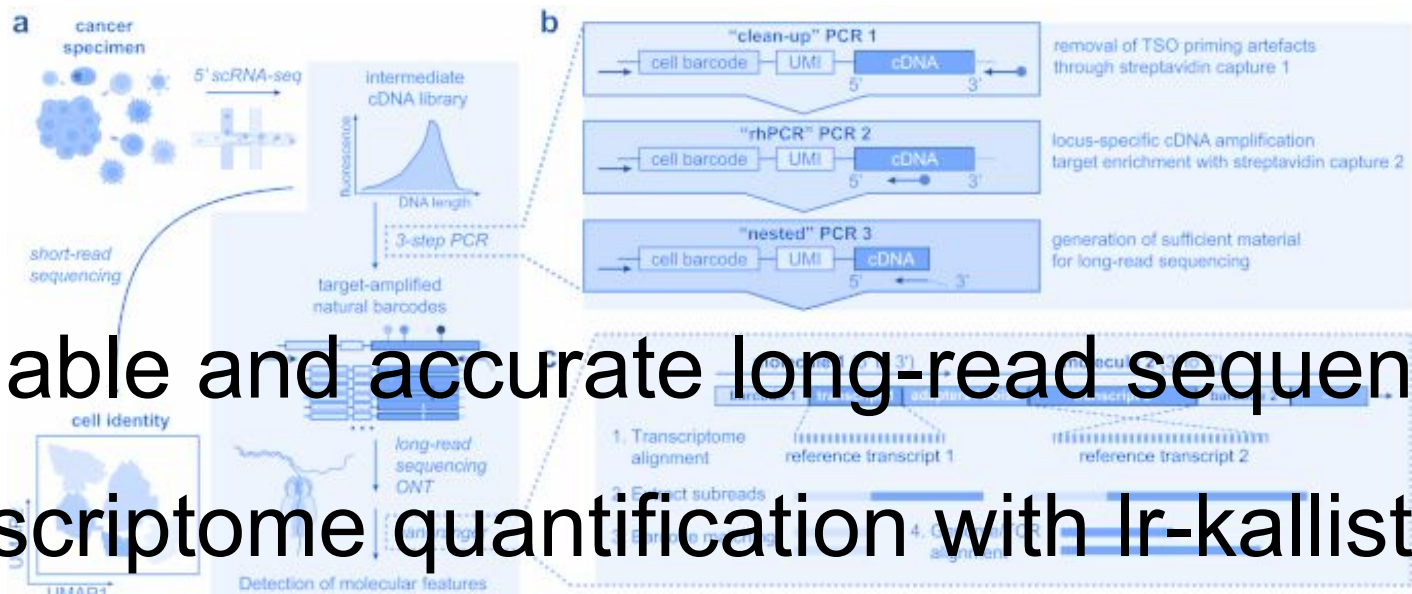
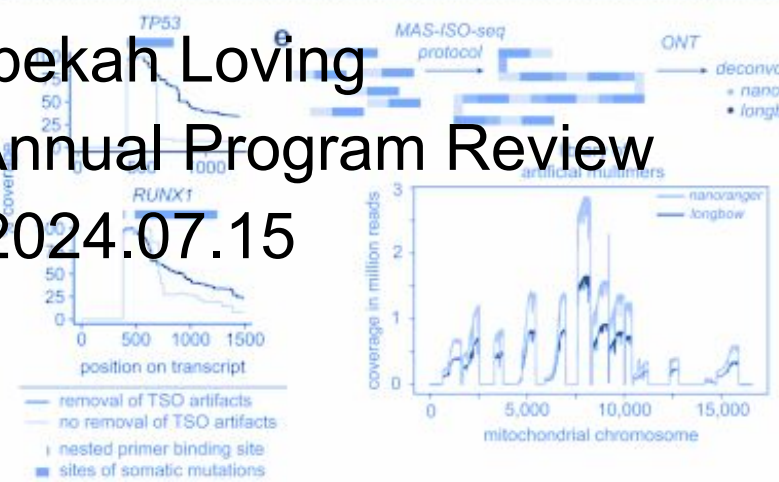


Scalable and accurate long-read sequencing transcriptome quantification with Ir-kallisto

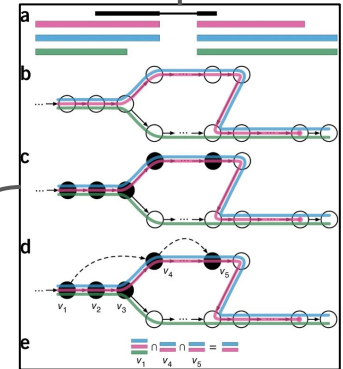
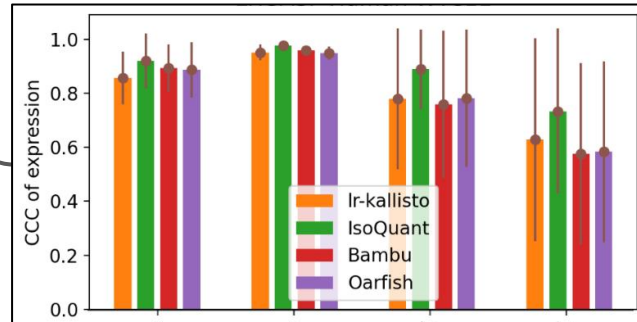
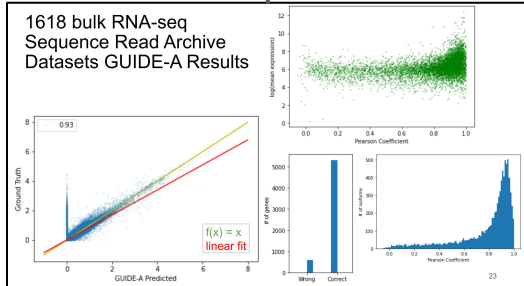
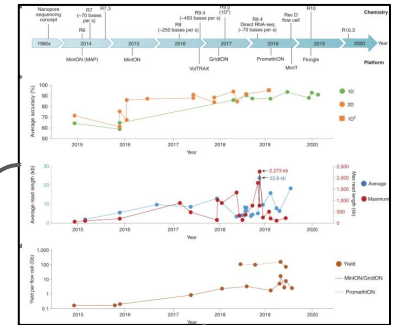
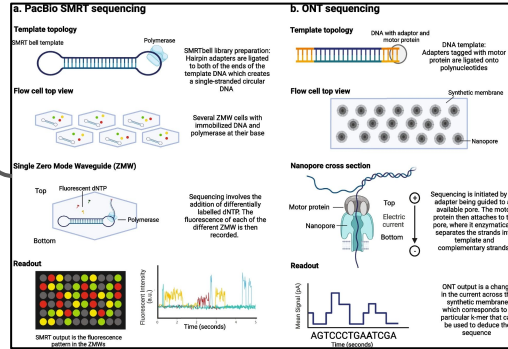
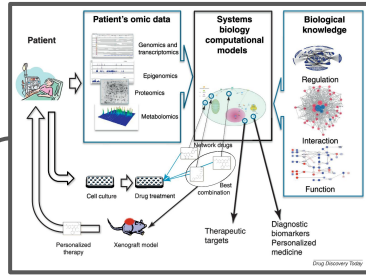
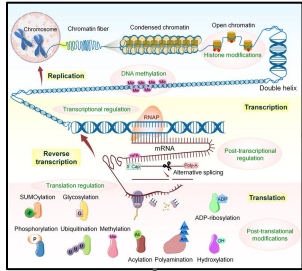


Rebekah Loving
 DOE CSGF Annual Program Review

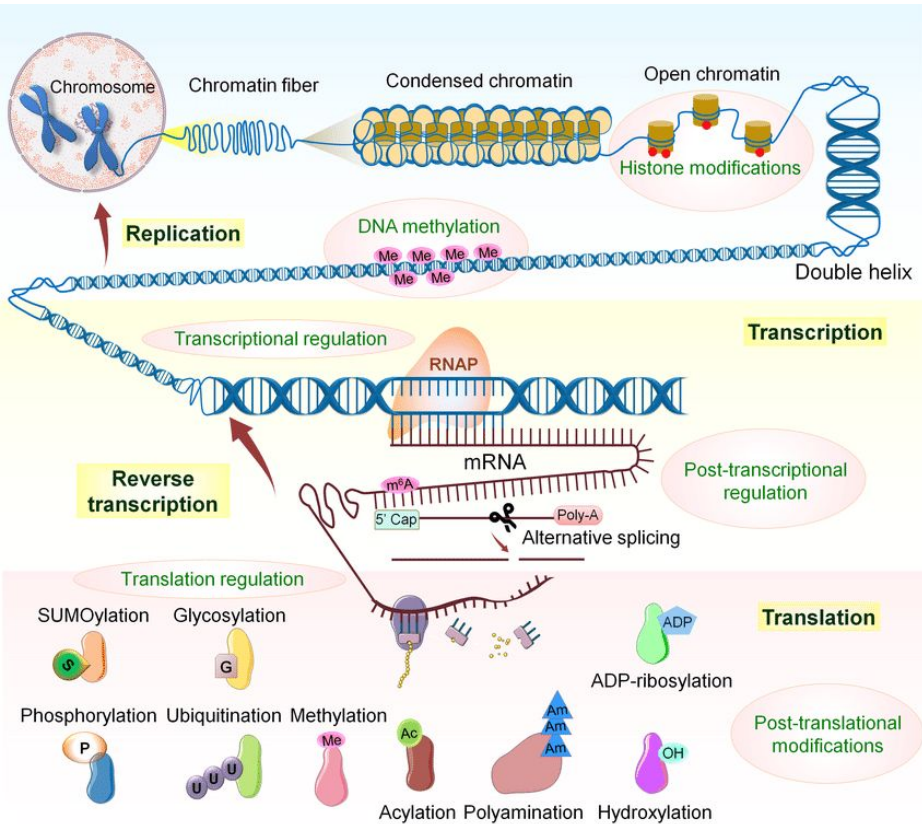
2024.07.15



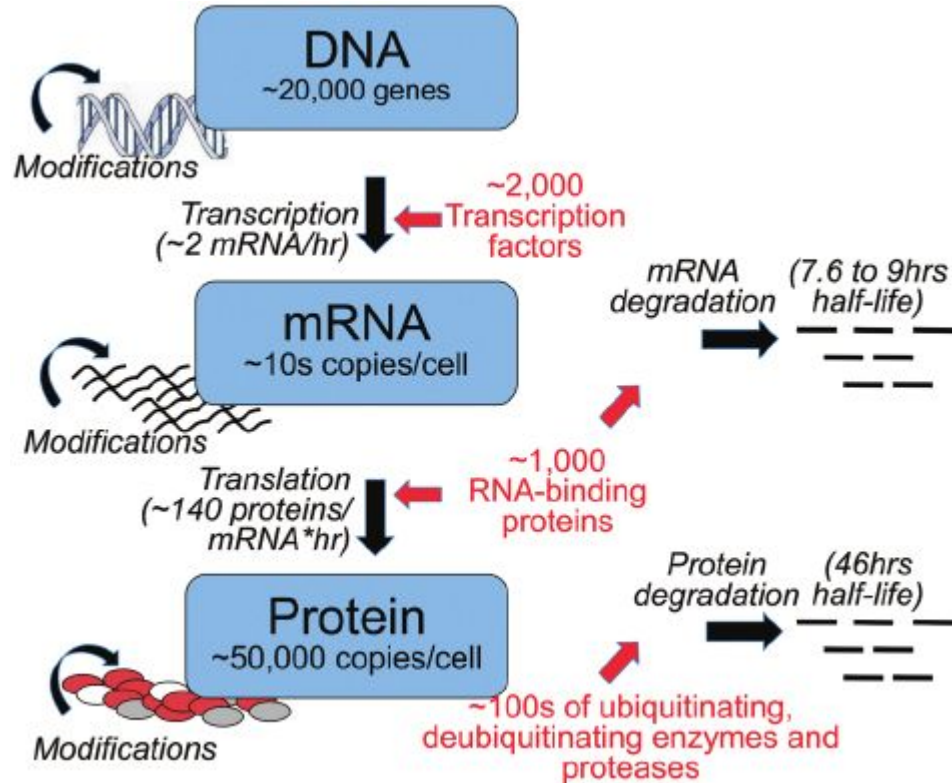
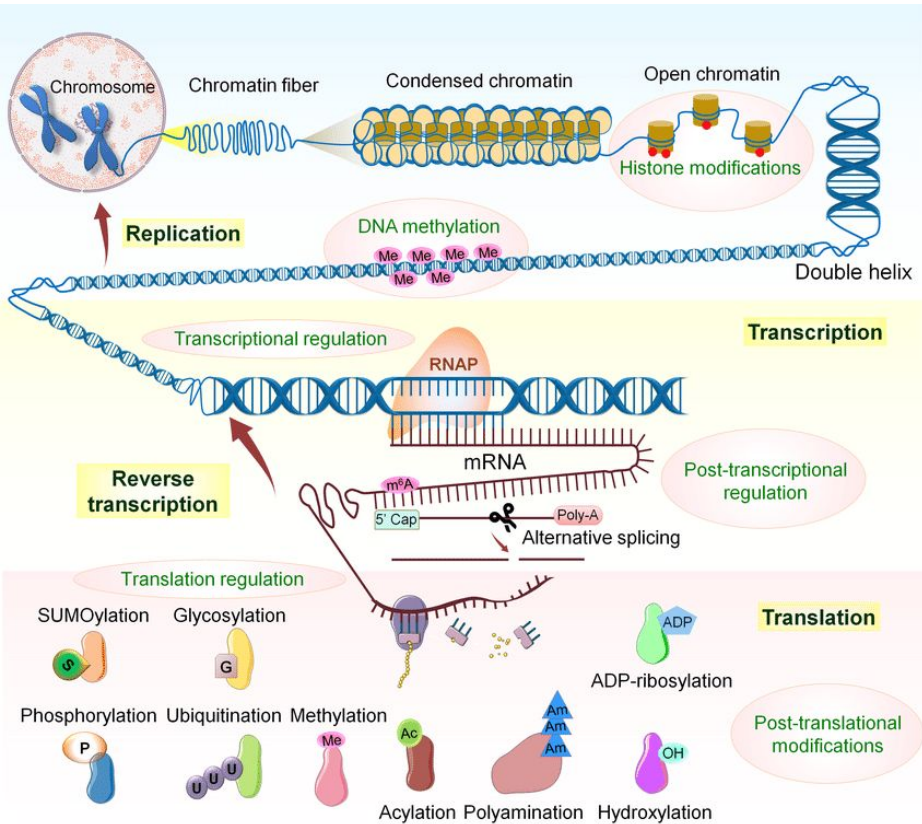
Roadmap



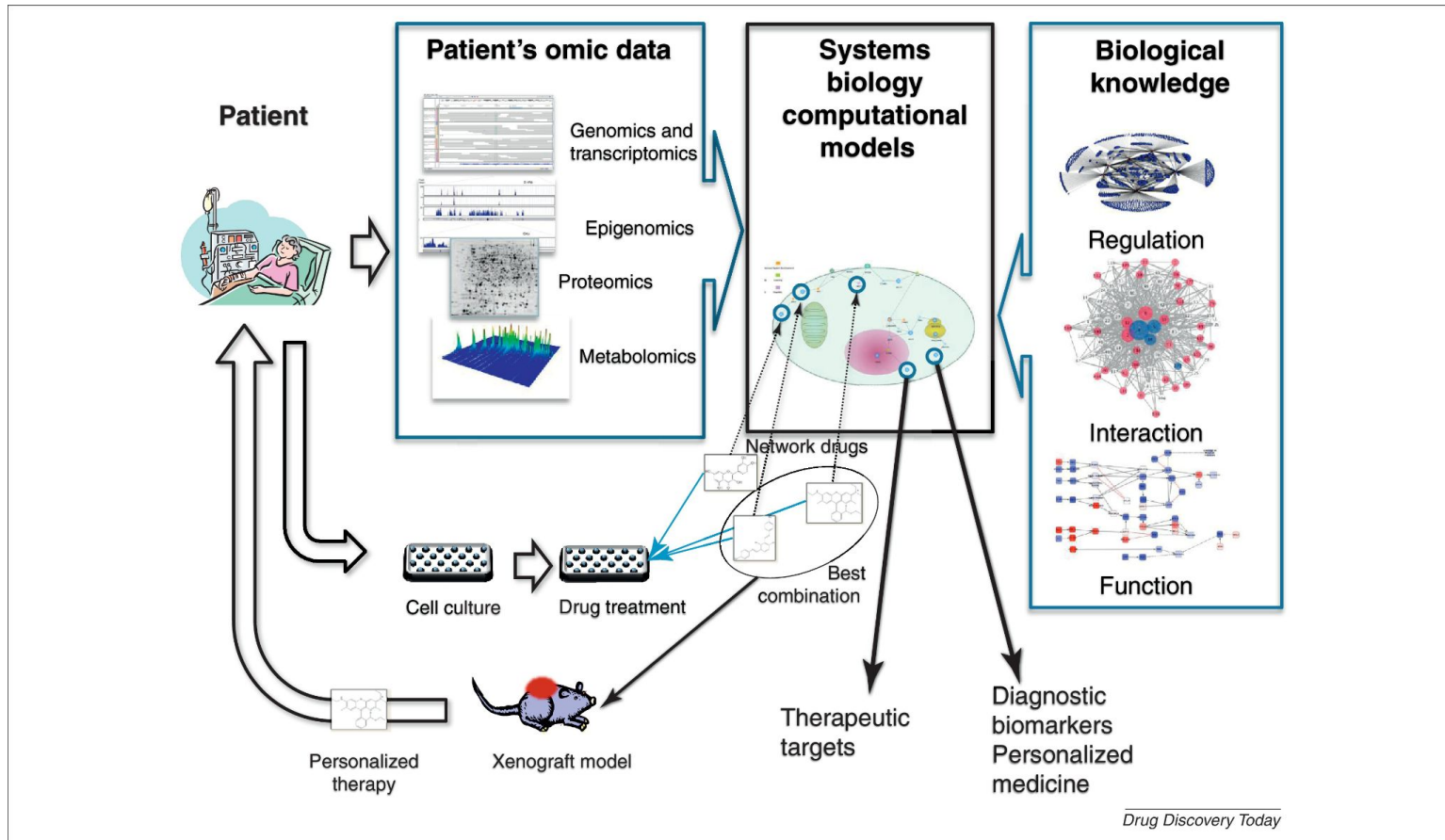
Central Dogma



Central Dogma



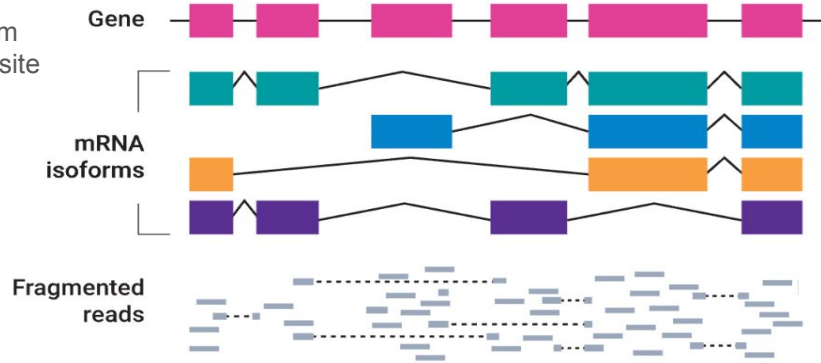
Studying mRNA impact on drug discovery



SHORT-READ SEQUENCING

LONG-READ SEQUENCING

Adapted from
PacBio website



**Identifying transcripts
is an assembly problem**



Partial view of isoform repertoire



Theoretically, no assembly required



Complete view of isoform repertoire

Long reads allow you to identify:

- Isoform resolution
- Improved annotations
- Gene fusions

New long read technology reduces error rate to .1-1.5%

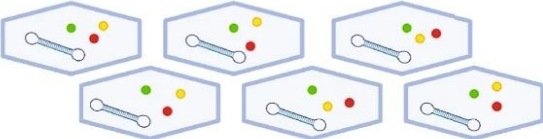
PacBio's innovation of Circular Consensus Sequencing produced an error rate .1-.5%

a. PacBio SMRT sequencing

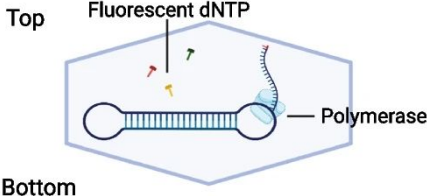
Template topology



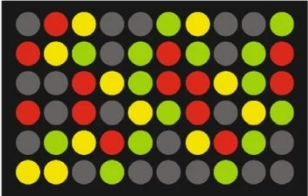
Flow cell top view



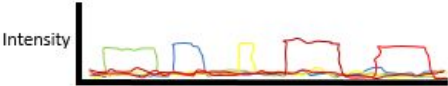
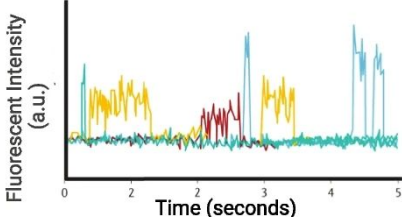
Single Zero Mode Waveguide (ZMW)



Readout



SMRT output is the fluorescence pattern in the ZMWs

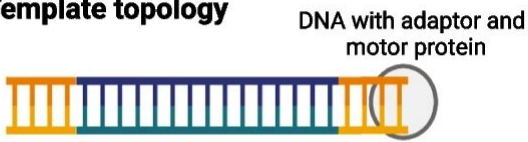


Adapted from PacBio website and Oehler, J.B., Wright, H., Stark, Z. et al. The application of long-read sequencing in clinical settings. *Hum Genomics* (2023).

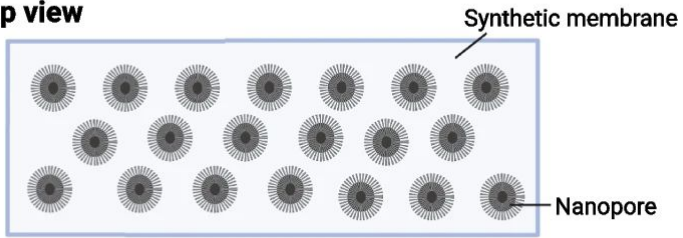
New long read technology reduces error rate to .1-1.5%

b. ONT sequencing

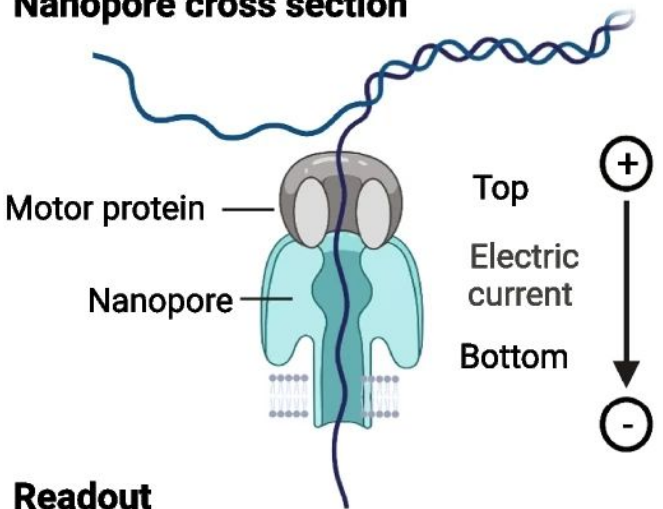
Template topology



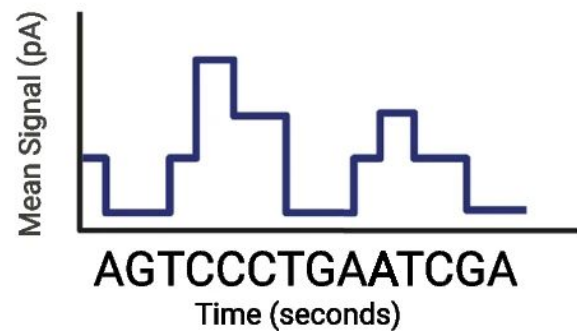
Flow cell top view



Nanopore cross section



Readout



Adapted from Oehler, J.B., Wright, H., Stark, Z. et al. The application of long-read sequencing in clinical settings. *Hum Genomics* (2023).

GUIDE-A Key Question

Does the 3' end alone hold sufficient information to predict isoforms' expression accurately, i.e. can the 3' end reveal key alternative splicing patterns?

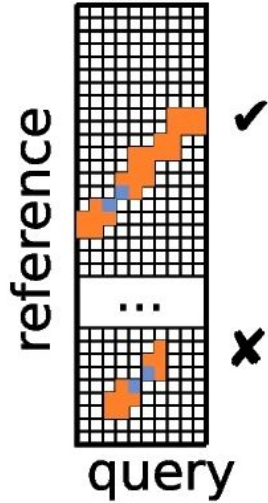
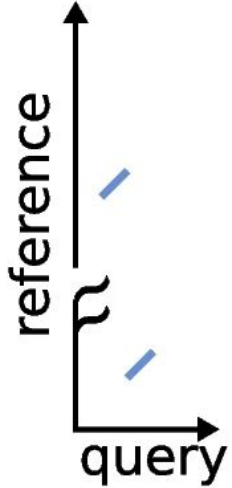
Challenges of RNA-seq mapping and quantification

Short read tools vs long read tools

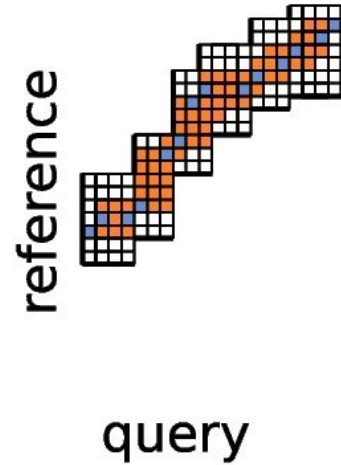
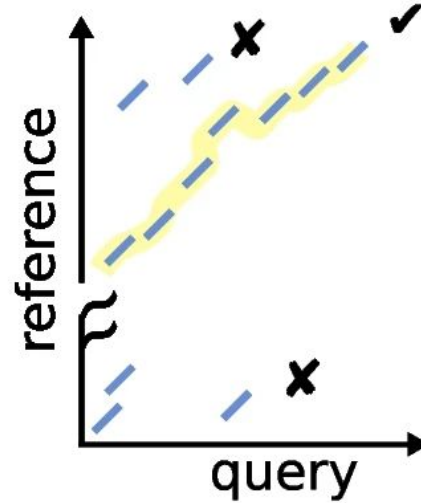
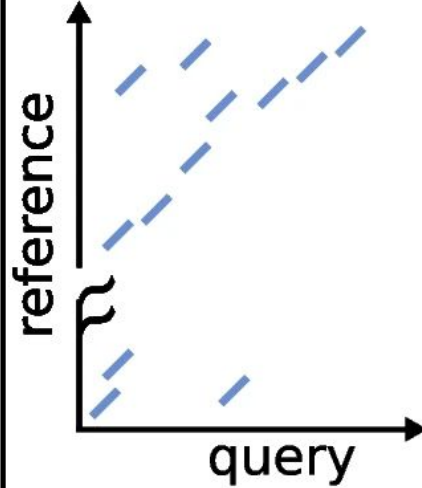
Challenges	kallisto	salmon	Oarfish	bambu	IsoQuant
Scalability and efficiency	Green	Green	Green	Green	Red
Accuracy	Green	Green	Red	Yellow	Green
Robustness (to error)	Grey	Grey	Red	Green	Yellow
Flexibility (to different assays)	Green	Red	Red	Yellow	Yellow
Alignment	Green	Green	minimap2	minimap2	minimap2

Alignment challenges

short-read mapping



long-read mapping



a) seeding

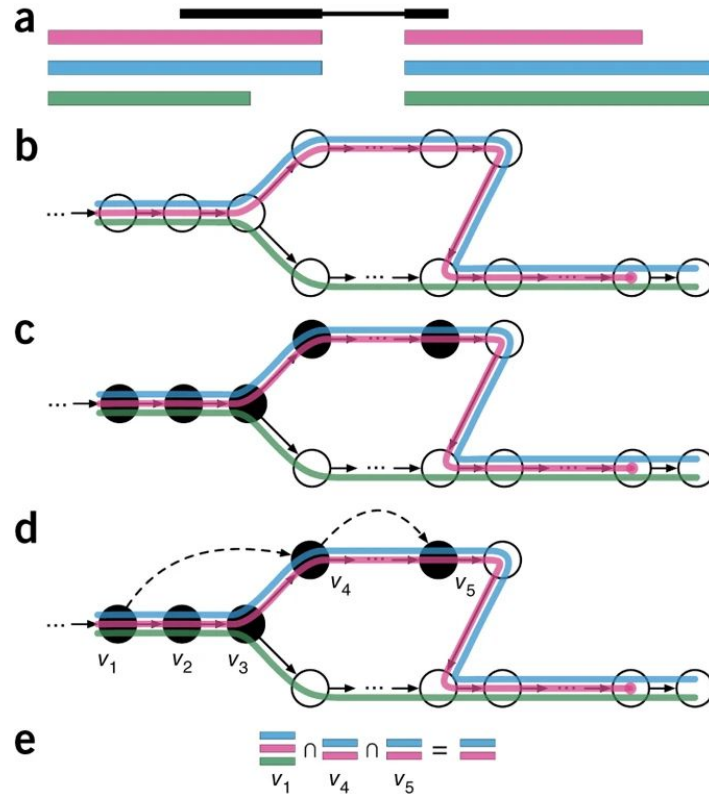
b) extending

a) seeding

b) chaining

c) extending

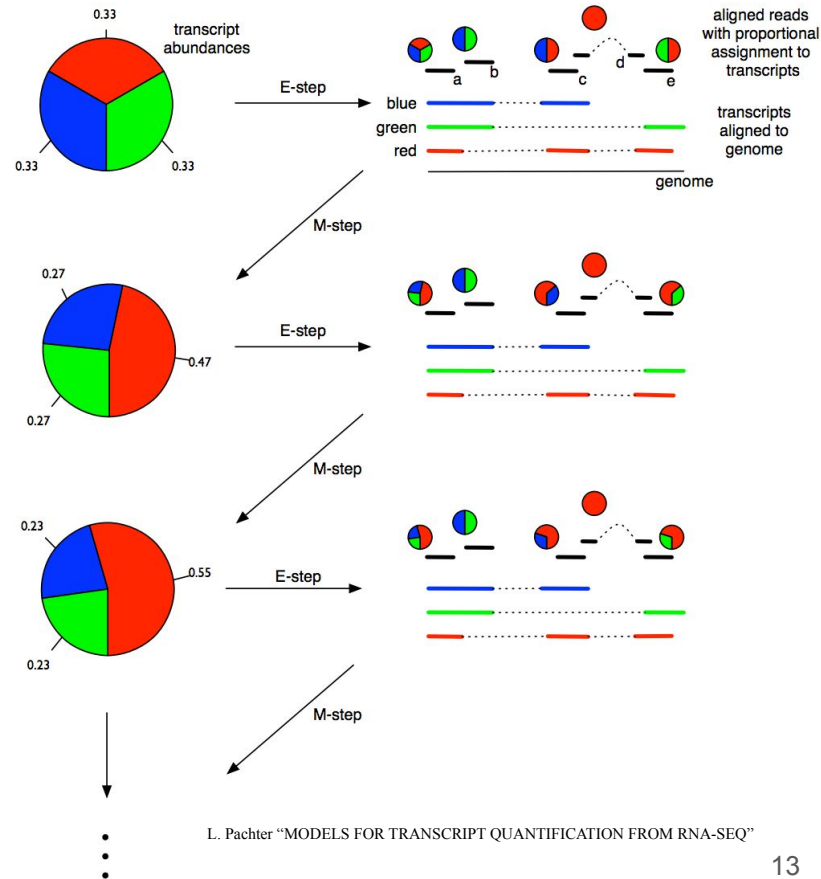
Ir-kallisto Background



Ir-kallisto Background

$$L(\alpha) \propto \prod_{f \in F} \sum_{t \in T} y_{f,t} \frac{\alpha_t}{l_t} = \prod_{e \in E} \left(\sum_{t \in e} \frac{\alpha_t}{l_t} \right)^{c_e}$$

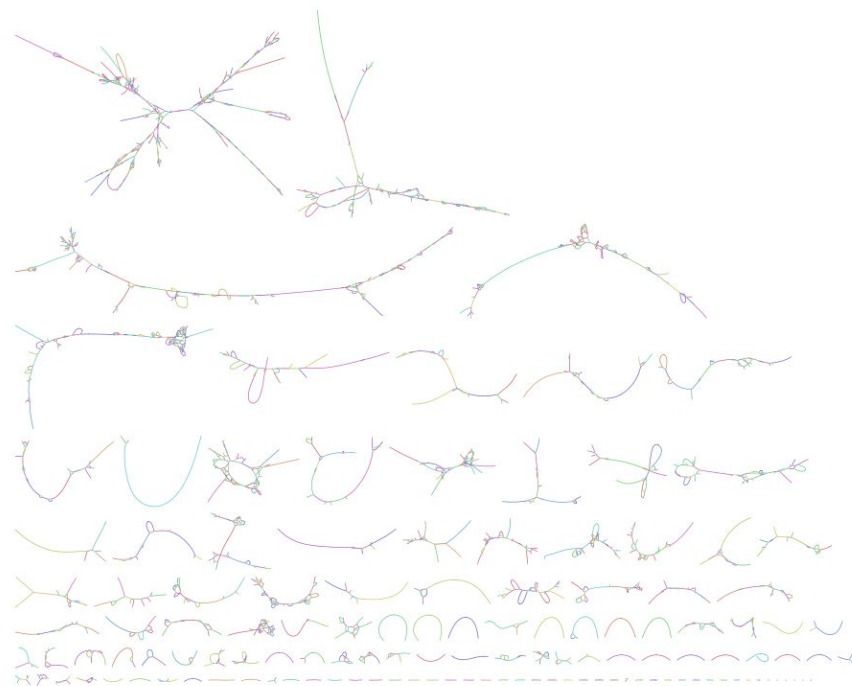
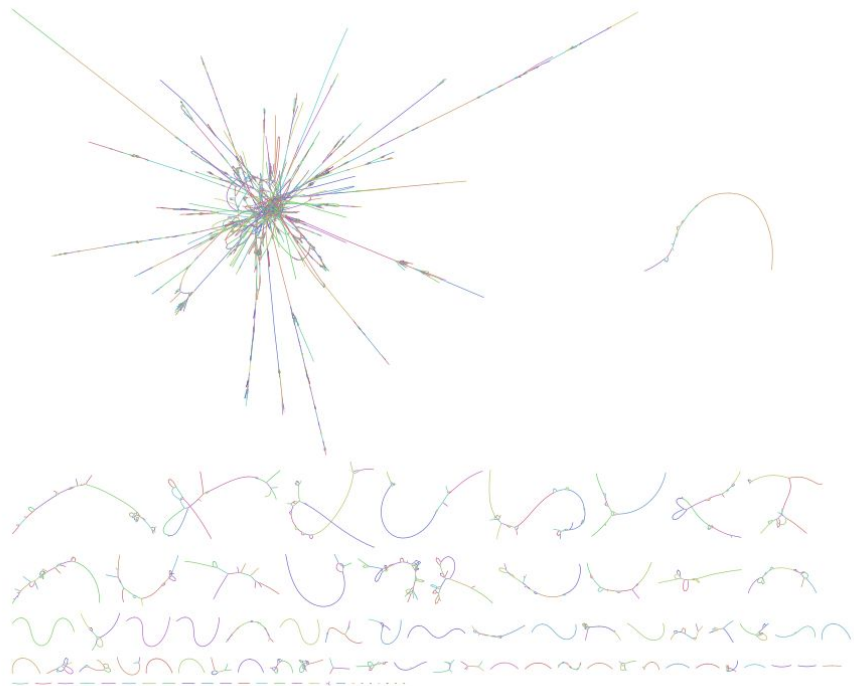
l_t is the (effective) length of transcript t



Effect of k -mer length on DBG in kallisto

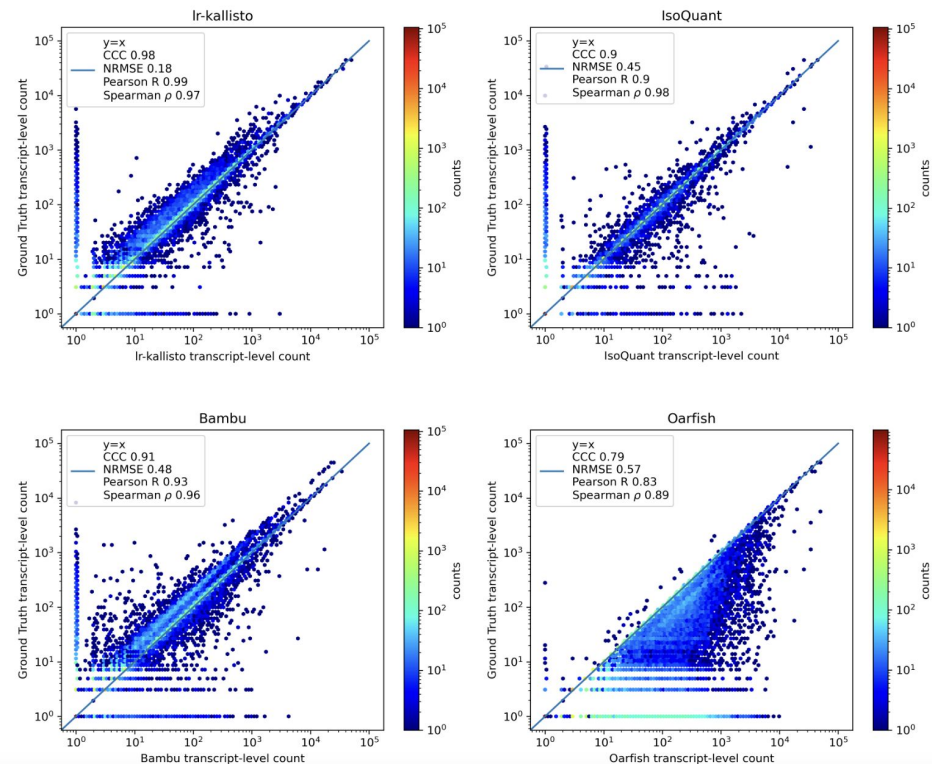
ii. Ex.: first 1000 transcripts, k -mer=31

k -mer=63

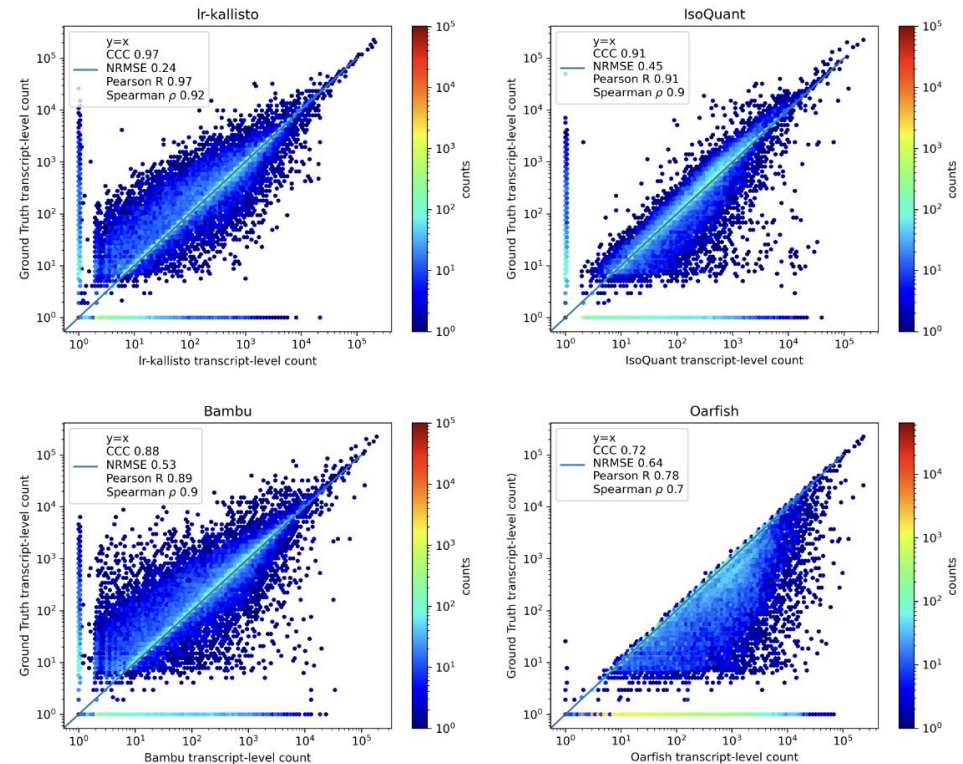


Simulation Benchmark

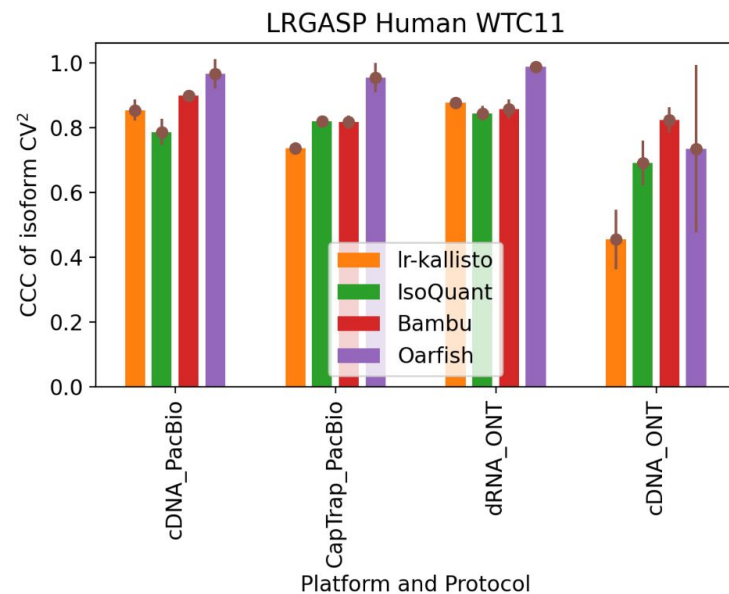
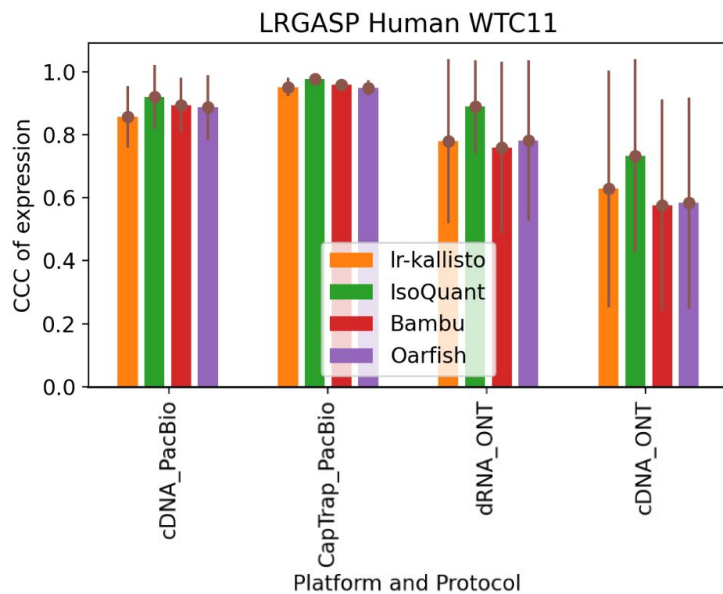
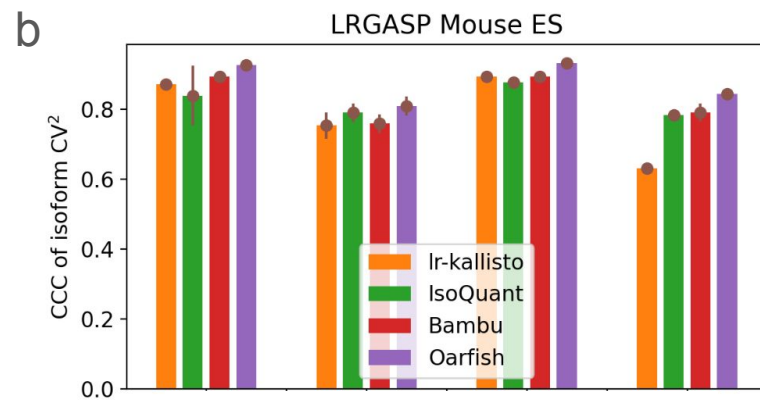
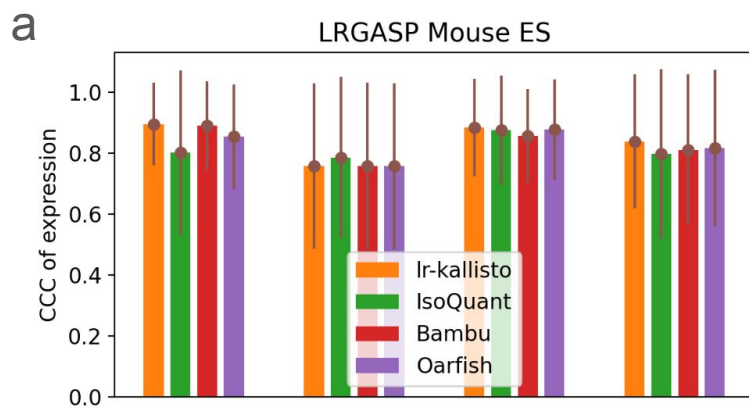
a PacBio IsoSeqSim Simulation



b ONT NanoSim R10.4 Simulation



LRGASP Reproducibility Benchmark



GUIDE-A: Gene Unambiguous Isoform Deduction Extraction - Algorithm

A machine learning approach for isoform prediction and gene/isoform network detection from 3' RNA-seq data

TCC = transcript compatibility counts
TPM = transcripts per million

GUIDE-A Problem Setup

Inputs: 3' end Transcript Compatibility Counts, **TCC**

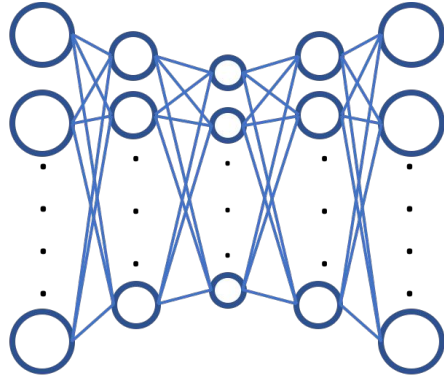
Outputs: full length Transcript Abundances, **TPM**

Find Mapping between **3' end** TCC and **full length** TPM

Mapping is complex; can we use a modified autoencoder to uncover the **correlations** and discover the mapping?

GUIDE-A

Architecture and Setup



Each layer is a **Linear Unit**.

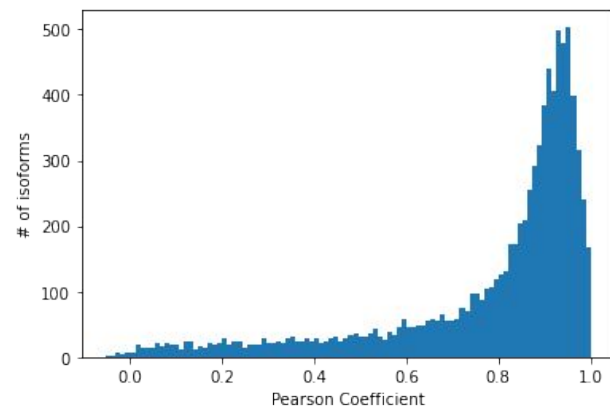
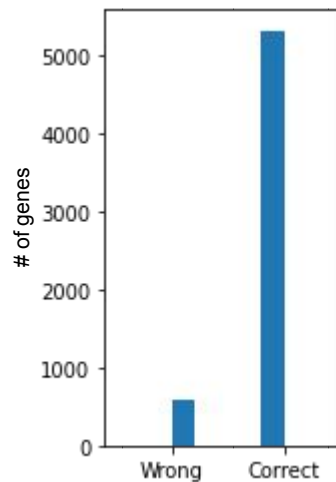
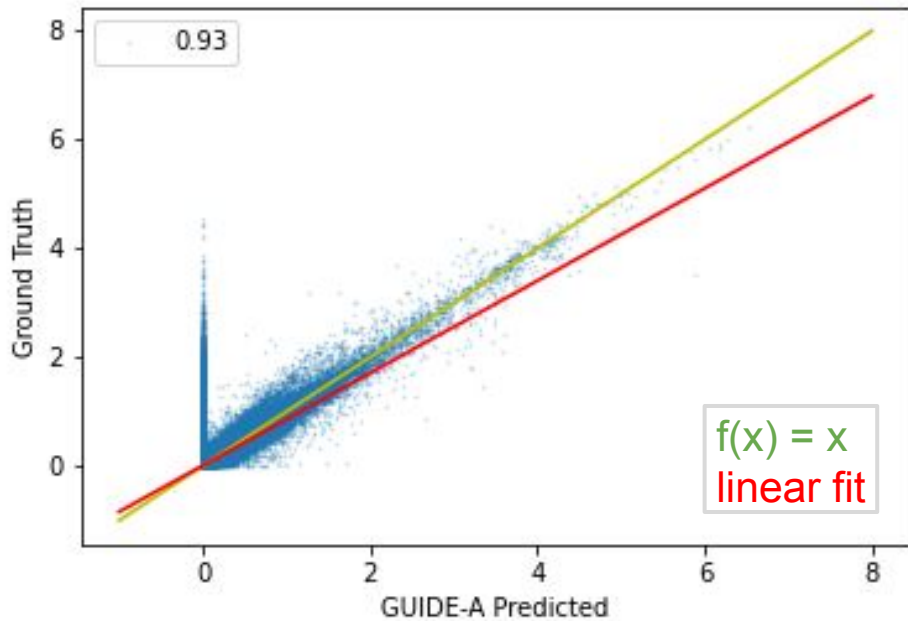
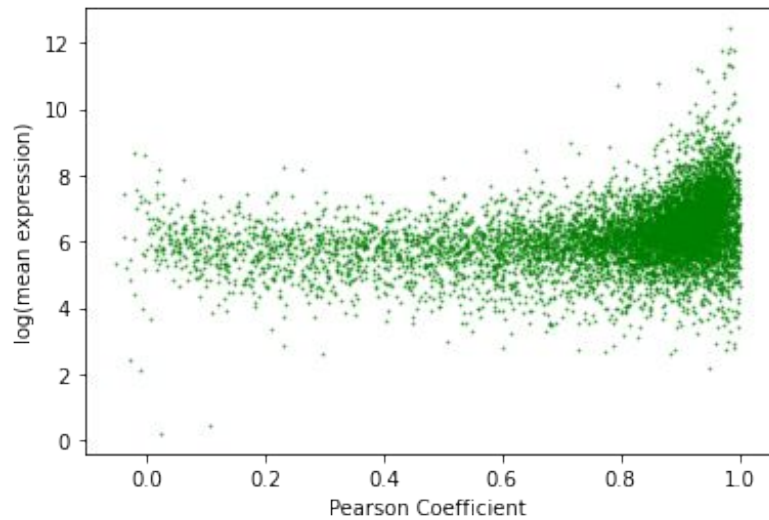
Rectified Linear Unit (ReLU) activation is used after each layer.

Batch normalization is used preceding the output layer.

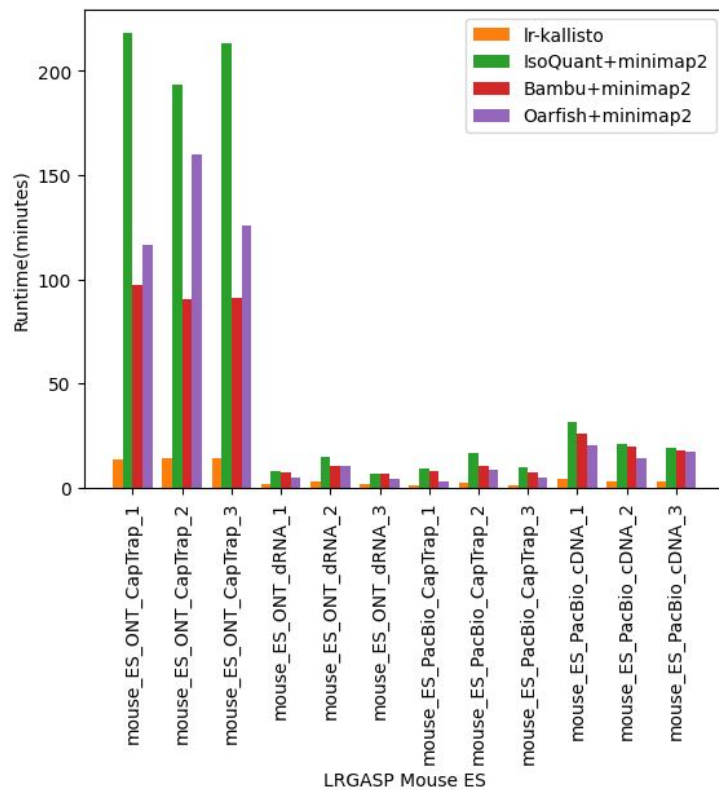
Loss is calculated via **Sum Squared Errors (SSE)**.

Training set is 80%, validation is 10%, and test is 10%, which was created by random sampling without replacement.

1618 bulk RNA-seq Sequence Read Archive Datasets GUIDE-A Results



Runtime Efficiency



Acknowledgements

Pachter Lab
Lior Pachter
Delaney
Laura
Tara
Kayla
Maria
Nikki
Ángel
Kristjan
Taleen
Gennady
Cat
Meichen
Lambda
Joe Rich

Mortazavi Lab
Ali Mortazavi
Dana Wyman
Fairlie Reese
Jaz Sakr
Liz Rebboah

Wold Lab
Barbara Wold
Diane Trout
Brian Williams

Committee Members:
Matt Thomson
Pietro Perona

Funding:
DOE CSGF
IGVF Consortium

