# KSBi-BIML 2021

Bioinformatics & Machine Learning (BIML)
Workshop for Life Scientists

생물정보학 & 머쉰러닝 워크샵(온라인)

Transcriptome Profiling by Nanopore Direct RNA Sequencing

장혜식







## Bioinformatics & Machine Learning for Life Scientists BIML-2021

안녕하십니까?

한국생명정보학회의 동계 워크샵인 BIML-2021을 2월 15부터 2월 19일까지 개최합니다. 생명정보학 분야의 융합이론 보급과 실무역량 강화를 위해 도입한 전문 교육 프로그램인 BIML 워크샵은 2015년에 시작하였으며 올해로 7차를 맞이하게 되었습니다. 유례가 없는 코로나 대유행으로 인해 올해의 BIML 워크숍은 온라인으로 준비했습니다. 생생한 현장 강의에서만 느낄 수 있는 강의자와 수강생 사이의 상호교감을 가질수 없다는 단점이 있지만, 온라인 강의의 여러 장점을 살려서 최근 생명정보학에서 주목받고 있는 거의 모든 분야를 망라한 강의를 준비했습니다. 또한 온라인 강의의한계를 극복하기 위해서 실시간 Q&A 세션 또한 마련했습니다.

BIML 워크샵은 전통적으로 크게 생명정보학과 AI, 두 개의 분야로 구성되어오고 있으며 올해 역시 유사한 방식을 채택했습니다. AI 분야는 Probabilistic Modeling, Dimensionality Reduction, SVM 등과 같은 전통적인 Machine Learning부터 Deep Learning을 이용한 신약개발 및 유전체 연구까지 다양한 내용을 다루고 있습니다. 생명정보학 분야로는, Proteomics, Chemoinformatics, Single Cell Genomics, Cancer Genomics, Network Biology, 3D Epigenomics, RNA Biology, Microbiome 등 거의 모든 분야가 포함되어 있습니다. 연사들은 각 분야 최고의 전문가들이라 자부합니다.

이번 BIML-2021을 준비하기까지 너무나 많은 수고를 해주신 BIML-2021 운영위원회의 김태민 교수님, 류성호 교수님, 남진우 교수님, 백대현 교수님께 커다란 감사를 드립니다. 또한 재정적 도움을 주신, 김선 교수님 (Al-based Drug Discovery), 류성호 교수님, 남진우 교수님께 감사를 표시하고 싶습니다. 마지막으로 부족한 시간에도 불구하고 강의 부탁을 흔쾌히 허락하시고 훌륭한 강의자료를 만드는데 노력하셨을 뿐만아니라 실시간 온라인 Q&A 세션까지 참여해 수고해 주시는 모든 연사분들께 깊이감사드립니다.

2021년 2월

한국생명정보학회장 김동섭

#### 강의개요

#### Transcriptome Profiling by Nanopore Direct RNA Sequencing

짧은 리드를 이용한 RNA-seq은 RNA 발현을 살펴볼때 저렴한 가격에 높은 해상도를 얻을 수 있어 널리 활용되고 있다. 다만, 상보DNA로 변환하는 과정에서 RNA수식 정보가 사라지고, 짧은 리드로 인해 멀리 떨어진 스플라이싱 등의 정보를 얻는 데한계가 있다. 나노포어 시퀀싱은 DNA 중합효소를 전혀 사용하지 않는 새로운 RNA 시퀀싱 기법을 제공한다. 직접 RNA 시퀀싱(Direct RNA Sequencing; DRS)은 실험기법이 단순하고 결과가 수시간 안에 빠르게 나올 뿐만 아니라, 긴 리드에서 RNA 수식을 살펴볼 수 있는 기회를 제공하고 있어 기존 방법으로 해결 불가능했던 문제들에 활용될 수 있는 잠재력이 있다.

본 강의에서는 나노포어 시퀀싱과 DRS이 나오게 된 과정, 작동원리, 데이터 특성, 사용상 주의점들을 익히고, 전사체 분석에 활용할 수 있는 방법들을 소개한다. 아직 표준화된 도구들이 등장하지 않고 다양한 다른 접근법과 프로그램이 혼용되어 사용되는 DRS의 특성을 고려해서 응용 방법들의 사용되고 있는 프로그램들이 기반으로 하고 있는 생화학적, 알고리즘적 원리들에 대해 간략히 설명한다. 이를 통해, DRS를 연구에 활용하고자 할 때 당시에 새로 나오게 될 도구들의 특성을 쉽게 파악하고 익힐 수 있도록 하는 것을 목표로 한다.

강의는 다음의 내용을 포함한다:

- 나노포어 시퀀싱의 원리와 역사
- 나노포어 시퀀싱 실험과 데이터의 특성
- 직접 RNA 시퀀싱의 활용 용도별 주요 도구 및 그 원리
- 나노포어 시퀀싱 신호 처리법 개요

#### \* 교육생준비물:

Zoom 수강 가능 노트북, 인터넷

\* 강의: 장혜식 교수 (서울대학교 생명과학부)

#### **Curriculum Vitae**

#### Speaker Name: Hyeshik Chang, Ph.D.



#### ▶ Personal Info

Name Hyeshik Chang
Title Assistant Professor
Affiliation Seoul National University

#### **▶** Contact Information

Address: 1 Gwanak-ro Gwanak-gu, Seoul, 08826

Email: hyeshik@snu.ac.kr Website: https://qbio.io

Research interest: High-throughput sequencing, post-transcriptional regulation, RNA-protein interaction

#### **Educational Experience**

1998–2007 B.S.E. in Information and Industrial Engineering, Yonsei University, Korea

2007–2009 M.S.E. in Bio and Brain Engineering, KAIST, Korea

2009–2014 Ph.D. in Biological Sciences, Seoul National University, Korea

#### **Professional Experience**

2001–2005 Software Developer, Solution Development Team, LinuxKorea, Inc.

2014–2019 Research Assistant Professor, IBS Center for RNA Research, Seoul National University

2018– Research Fellow, Center for RNA Research, Institute for Basic Science

2019– Assistant Professor, School of Biological Sciences, Seoul National University

#### **Selected Publications (5 maximum)**

- 1. D. Kim, J.-Y. Lee, J.-S. Yang, J. W. Kim, V. N. Kim, and H. Chang. (2020) "The Architecture of SARS-CoV-2 Transcriptome." Cell, 181(4):914–921.
- 2. H. Chang<sup>1</sup>, J. Yeo<sup>1</sup>, J.-G. Kim, H. Kim, M. Lee, J. Lim, H. H. Kim, J. Ohk, H.-Y. Jeon, H. Lee, H. Jung, K.-W. Kim, and V. N. Kim. (2018) "Terminal uridylyltransferases execute programmed clearance of maternal transcriptome in vertebrate embryos." Molecular Cell, 70:72–82.e7.
- 3. J. Lim<sup>1</sup>, M. Ha<sup>1</sup>, H. Chang<sup>1</sup>, S. C. Kwon, D. K. Simanshu, D. J. Patel, and V. N. Kim. (2014) "Uridylation by TUT4 and TUT7 marks mRNA for degradation." Cell, 159(6):1365–1376.
- 4. H. Chang<sup>1</sup>, J. Lim<sup>1</sup>, M. Ha, and V. N. Kim. (2014) "TAIL-seq: genome-wide determination of poly(A) tail length and 3' end modifications." Molecular Cell, 53(6):1044–1052.
- 5. J. Cho<sup>1</sup>, H. Chang<sup>1</sup>, S. C. Kwon, B. Kim, Y. Kim, J. Choe, M. Ha, Y. K. Kim, and V. N. Kim. (2012) "LIN28A is a suppressor of ER-associated translation in embryonic stem cells." Cell, 151(4):765–777.

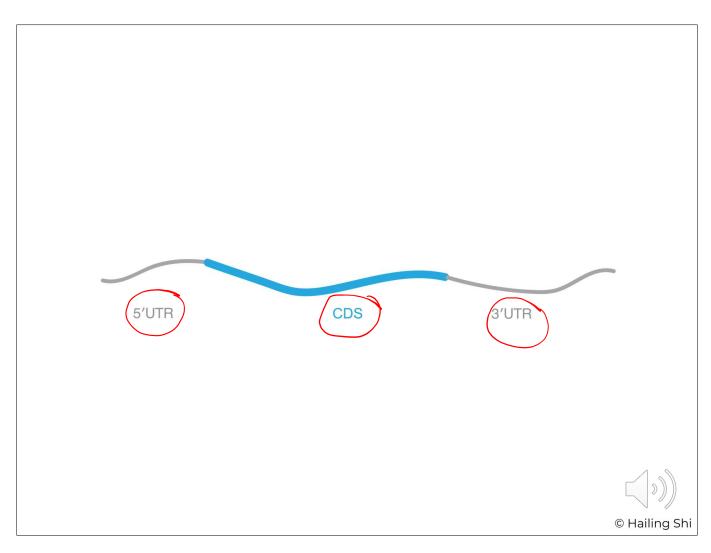
<sup>&</sup>lt;sup>1</sup> Co-first authors

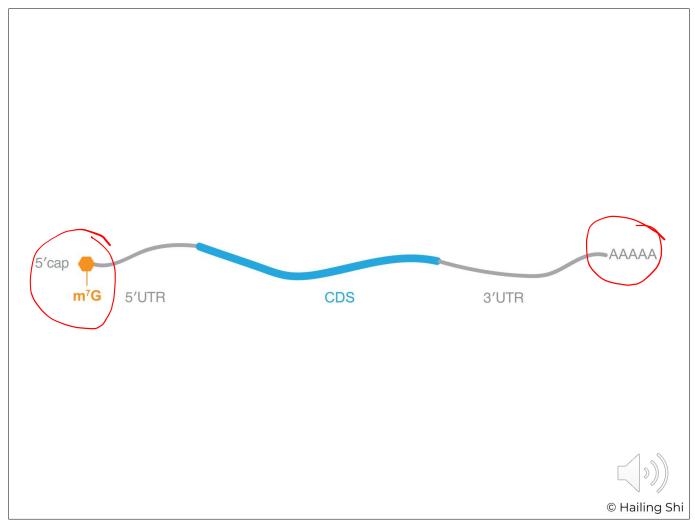


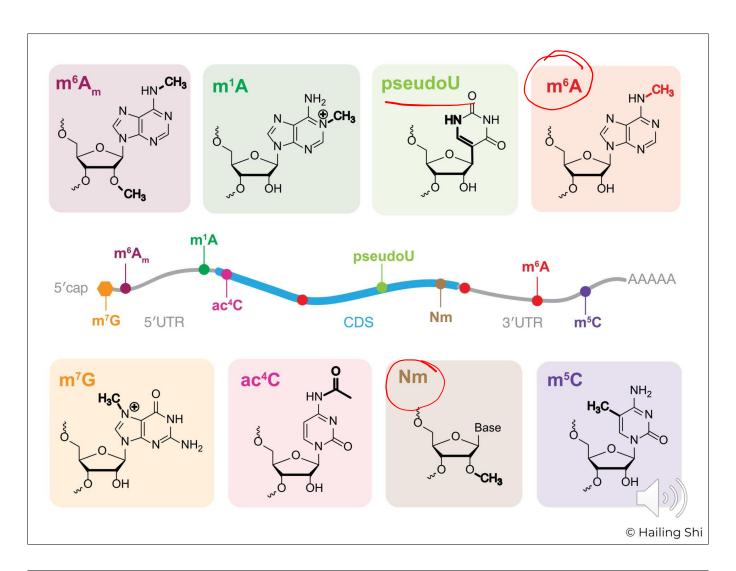
본 강의 자료는 한국생명정보학회가 주관하는 KSBi-BIML 2021 워크샵 온라인 수업을 목적으로 제작된것으로 해당 목적 이외의 다른 용도로 사용할 수 없음을 분명하게 알립니다. 수업 목적으로 배포 및 전송 받은 경우에도 이를 다른 사람과 공유하거나 복제, 배포, 전송할 수 없습니다.

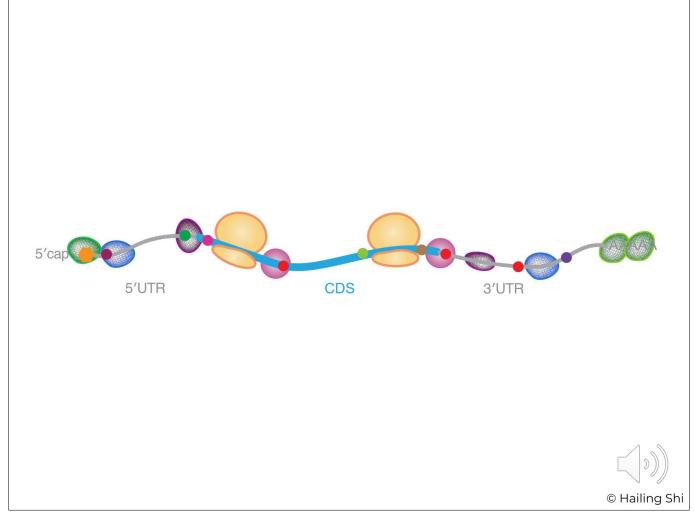
만약 이러한 사항을 위반할 경우 발생하는 모든 법적 책임은 전적으로 불법 행위자 본인에게 있음을 경고합니다.











Genomic DNA	Messenger RNA
Epigenome	Epitranscriptome (controversial)



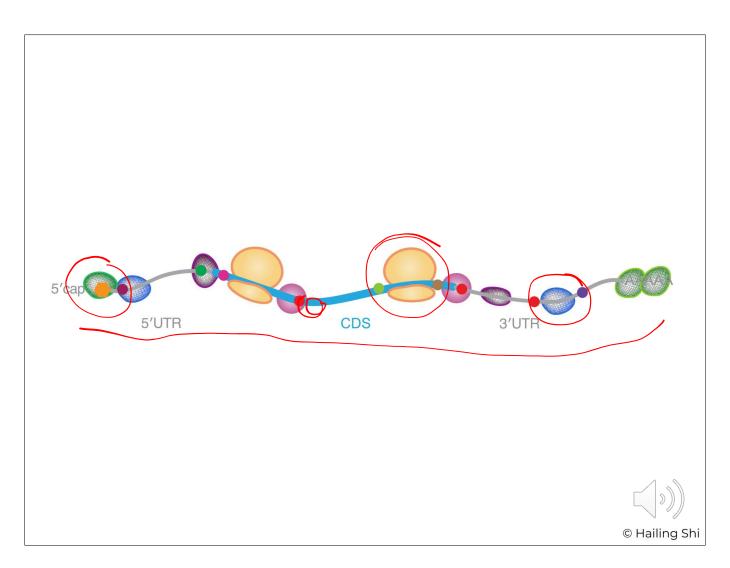
Genomic DNA	Messenger RNA
Epigenome	Epitranscriptome (controversial)
DNA Methylation	RNA Methylation RNA Editing

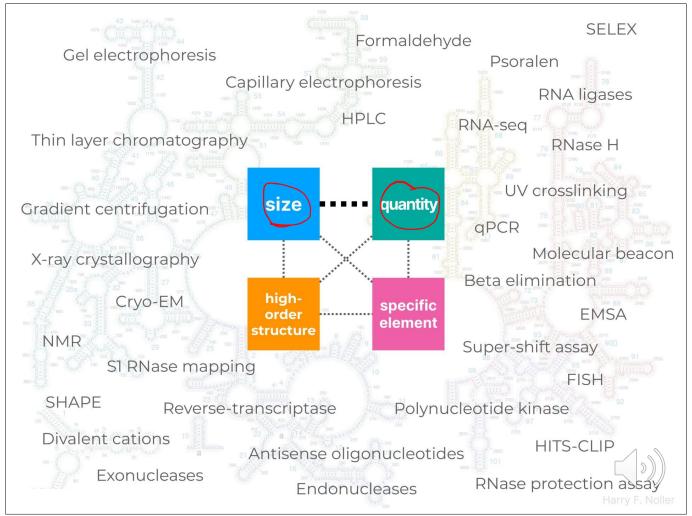


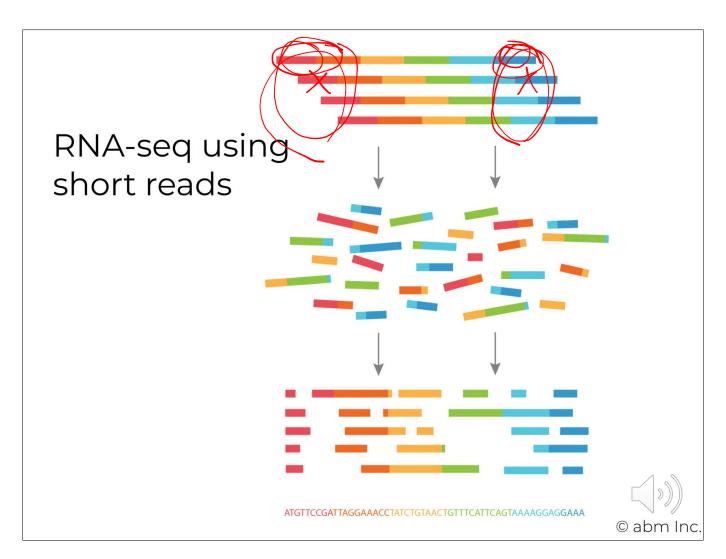
Genomic DNA	Messenger RNA
Epigenome	Epitranscriptome (controversial)
DNA Methylation	RNA Methylation RNA Editing
Histone Modification	Combinations of RBP Binding

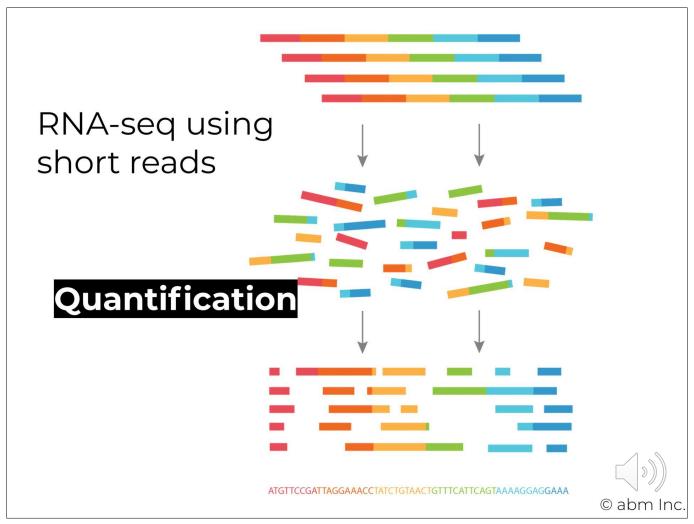


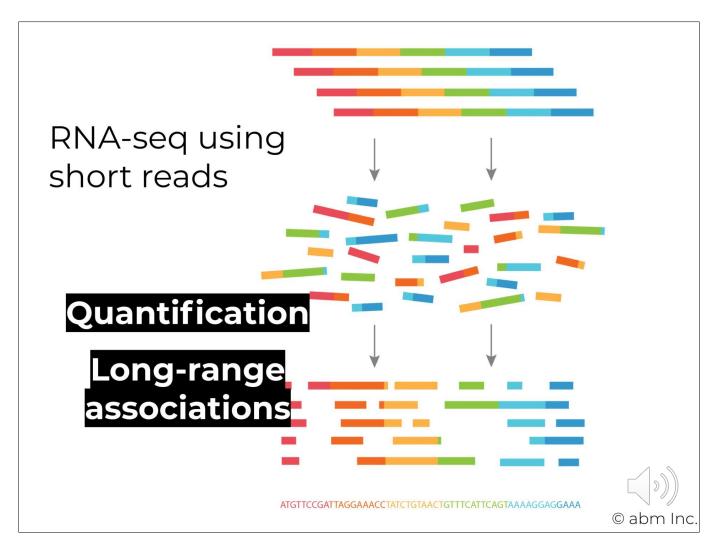
Genomic DNA	Messenger RNA
Epigenome	Epitranscriptome (controversial)
DNA Methylation	RNA Methylation RNA Editing
Histone Modification	Combinations of RBP Binding
Controls Transcription	Controls Splicing, Localization, Translation, and mRNA degradation

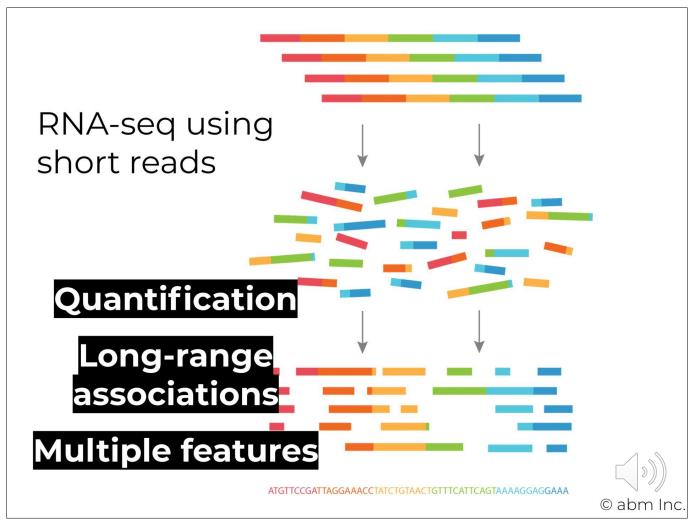




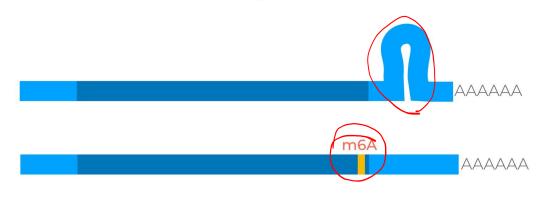






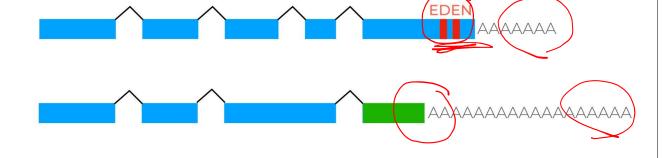


Does the **m6A** modifications near the stop codon affect the **secondary structure** of 3' UTR?

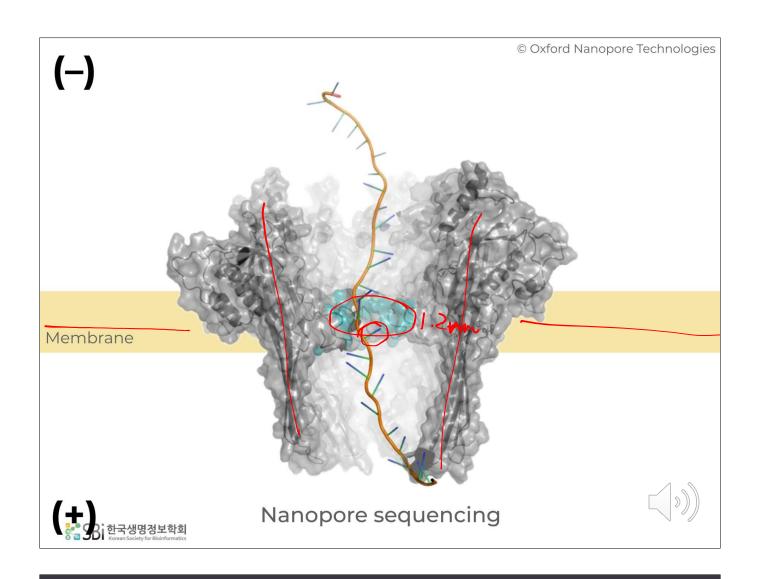




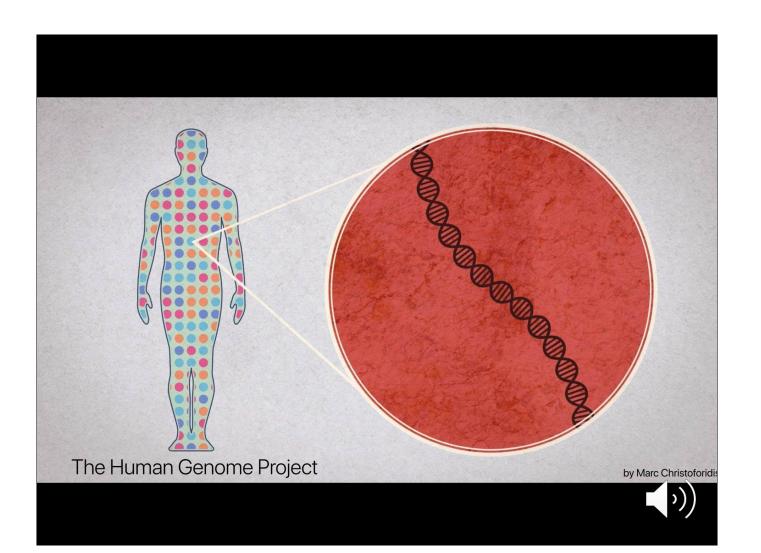
What **elements** make poly(A) lengths different for transcripts of the same gene?

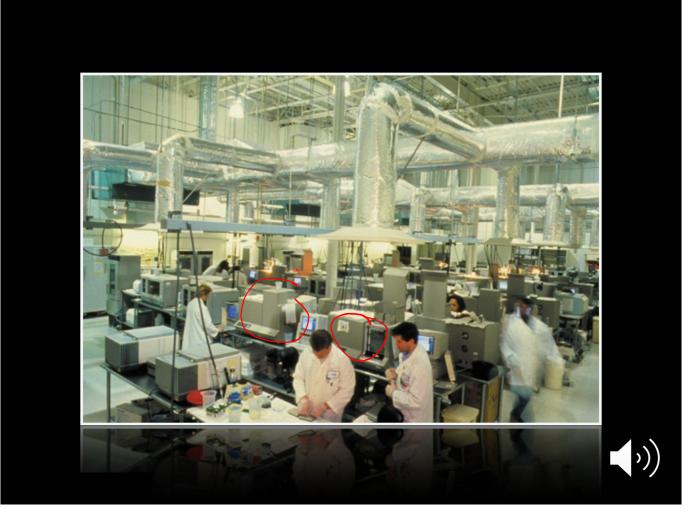


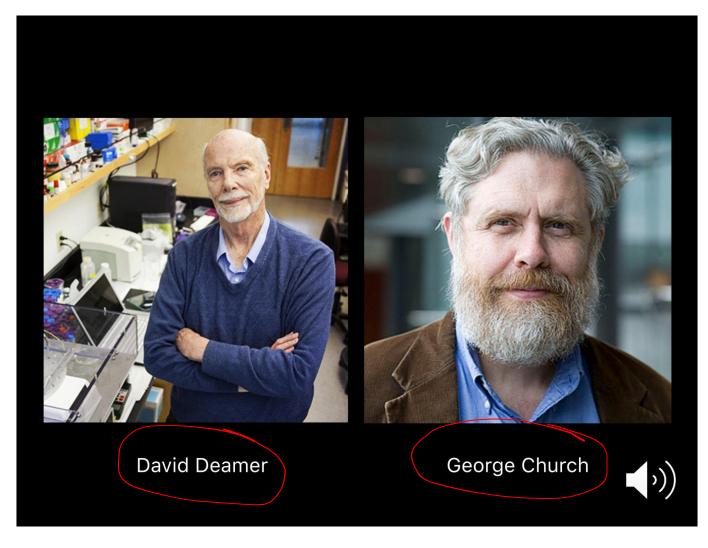


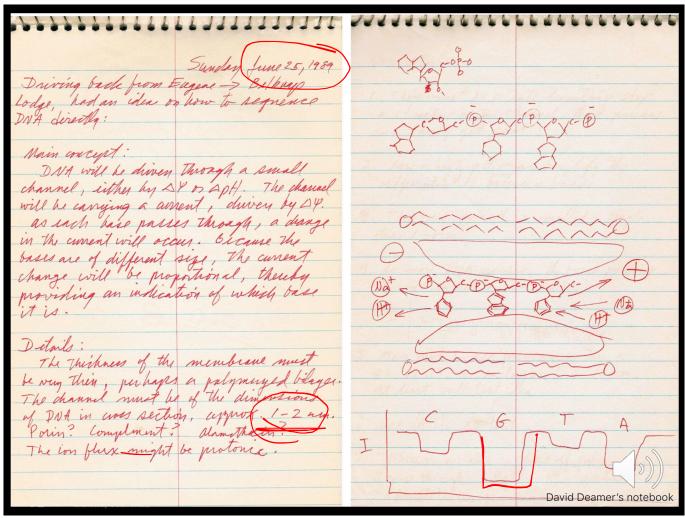


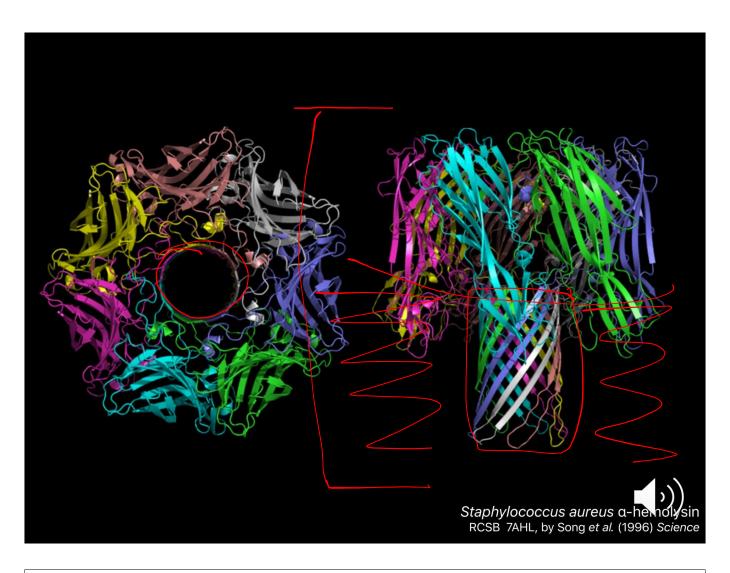


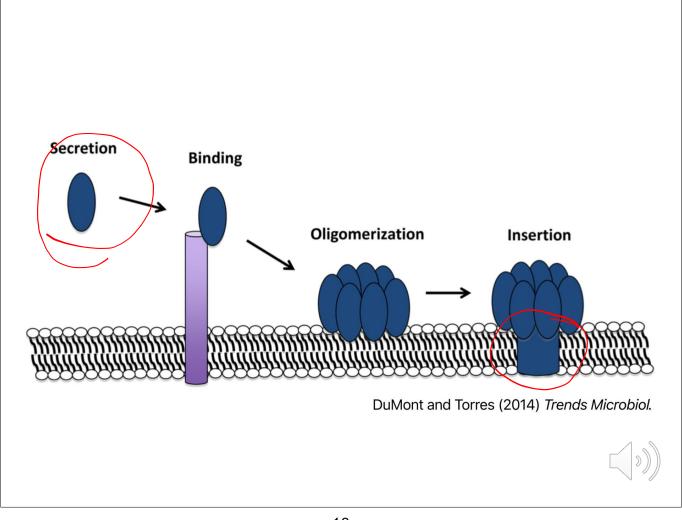






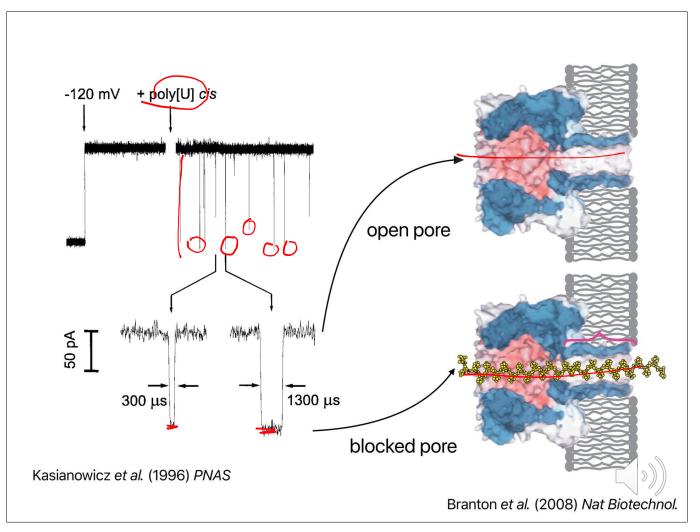


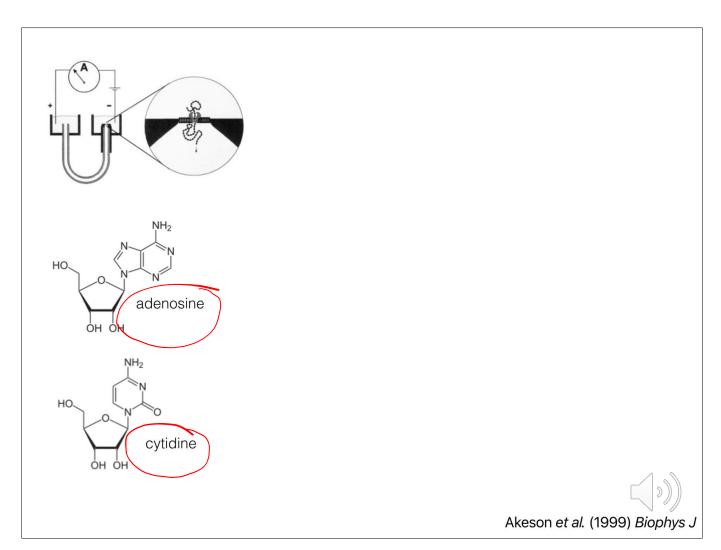


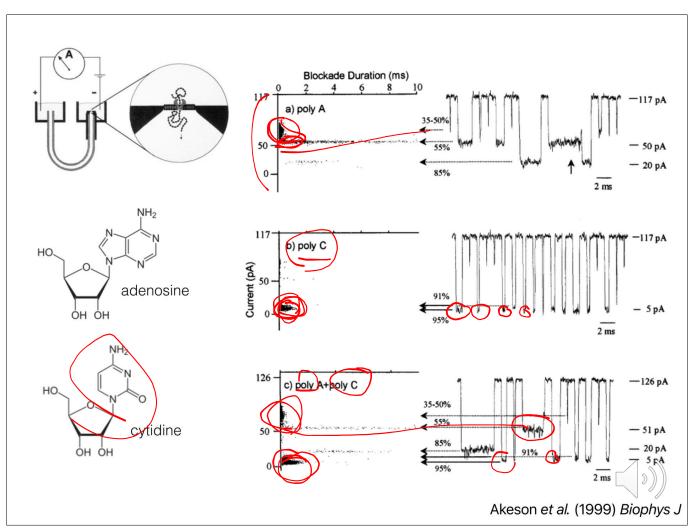


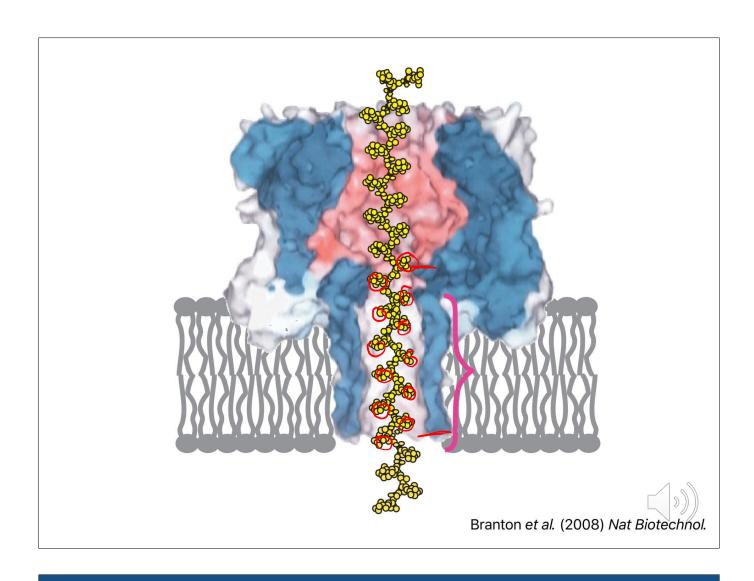
# What happens when a DNA passes through a nanopore?





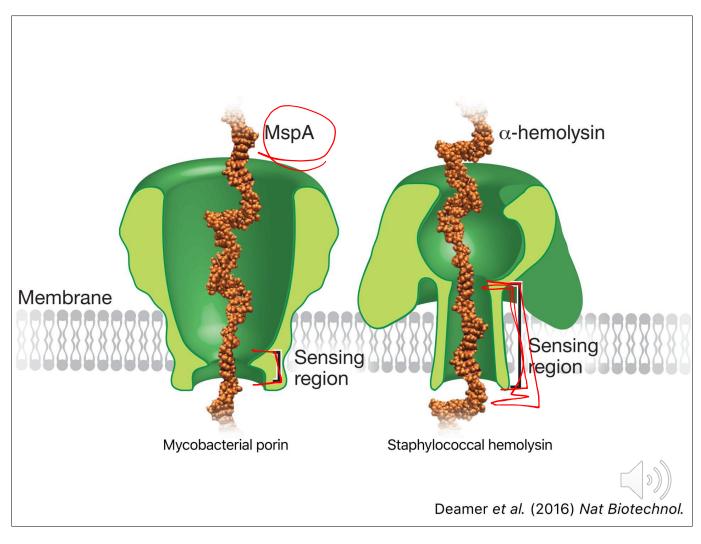


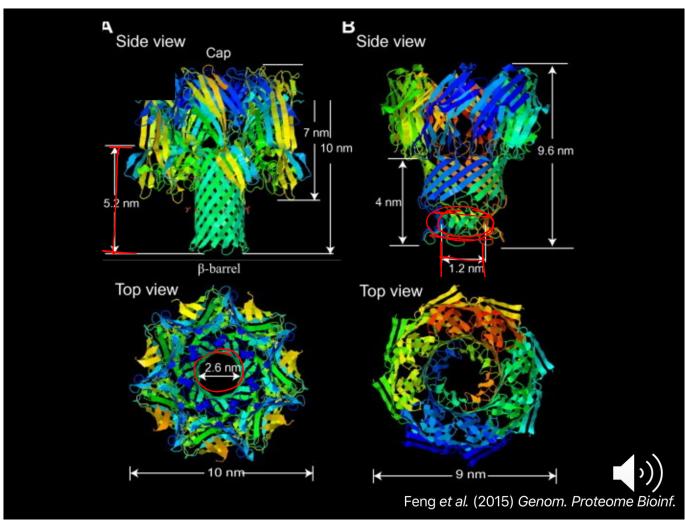


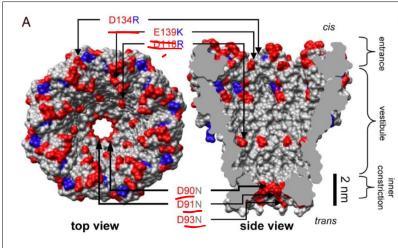


Could it be a bit more distinguishable?



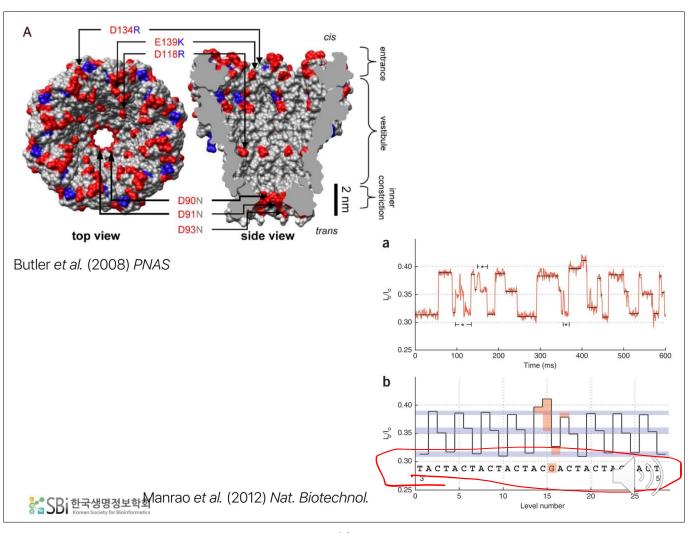


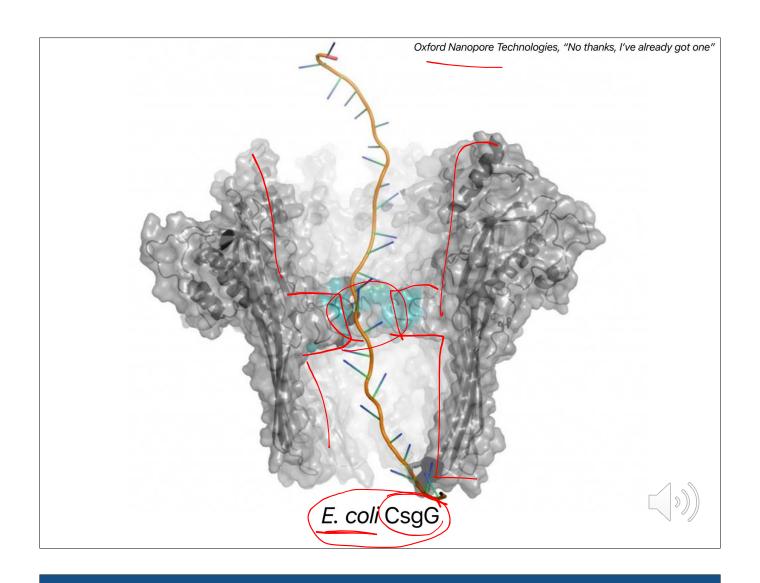




Butler et al. (2008) PNAS

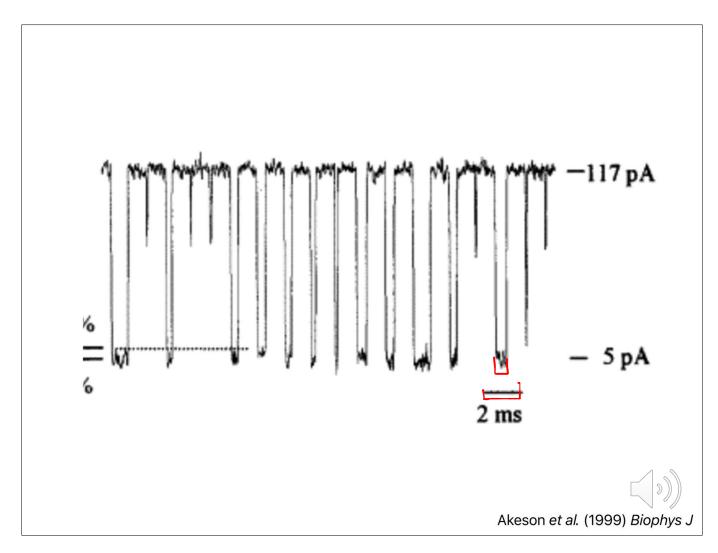


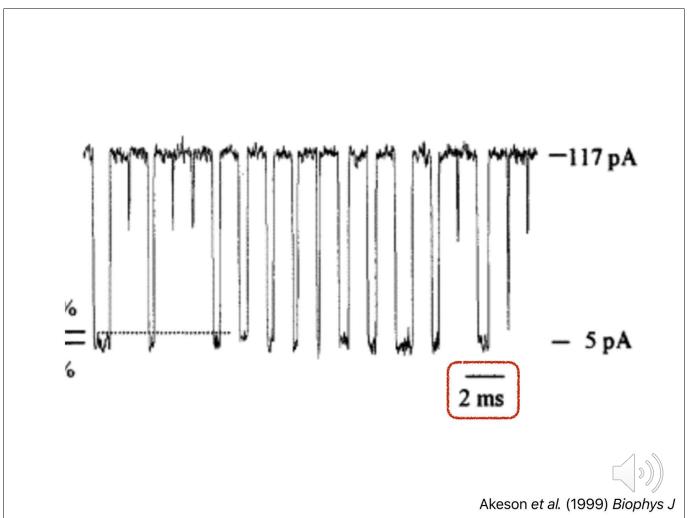


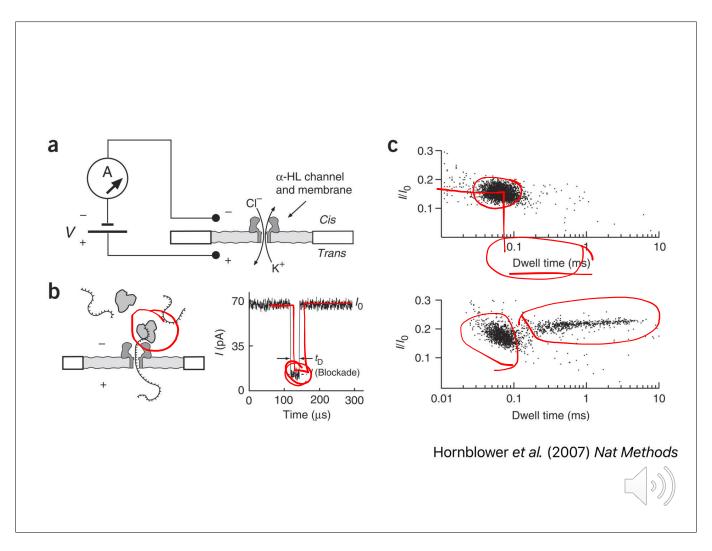


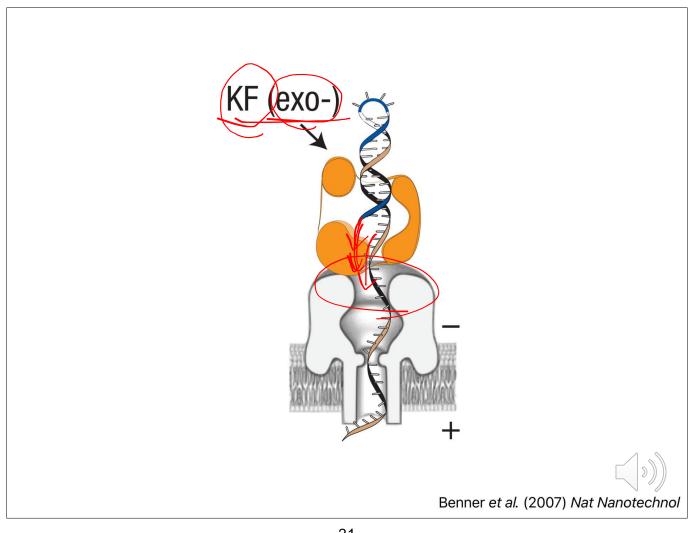
Could it be a bit more readable?

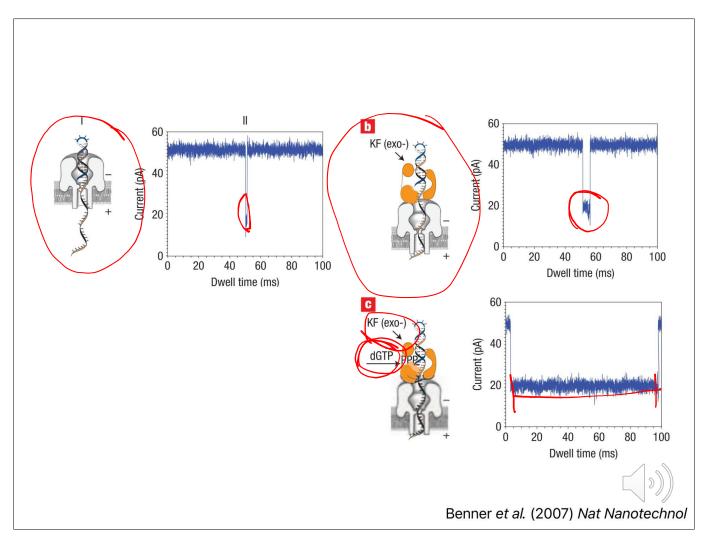


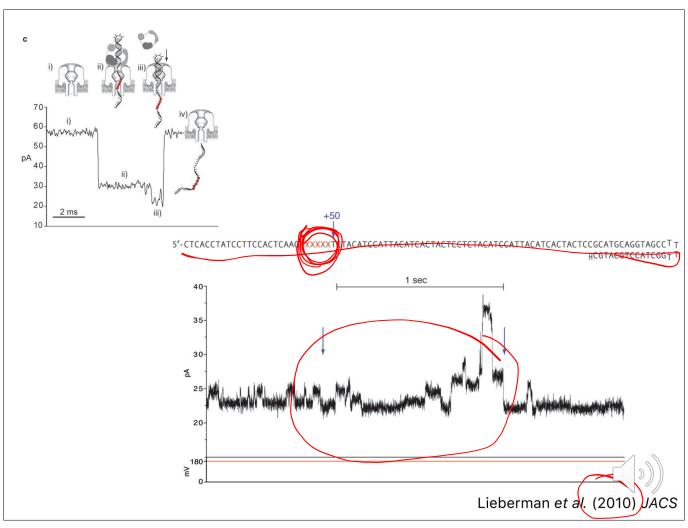






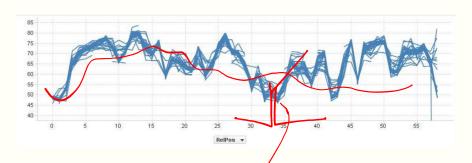






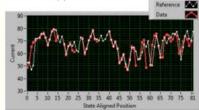
#### **Effect of Poor Movement on Base-calling**

- Missing states from problems with movement make base calling from single molecule data diffusione 2011
- Consensus plots can be generated from multiple single molecule reads



Base calls can be made from a consensus of single molecule reads by mapping to known model

Consensus mapped to static:



- Gaps in the sequence are evident even for known model
  - Need to make improvements to movement to drive progress

Oxford Nanopore Technologies, KeyGene Ser-

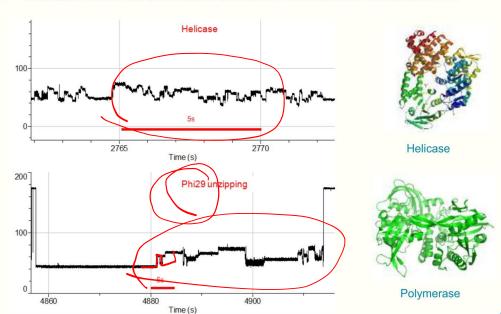
ar (2()16)

Confidential

#### Strand Sequencing – Movement July 2011

#### Comparison Between Phi29 and Helicase Data

- Problems with using Phi29 as a motor include:
  - · missed states, fast and slow movement regimes, backwards movement at low potential, pausing
- Initial work on helicases show that the distribution of movement is a lot more controlled



Oxford Nanopore Technologies, KeyGene Seminar (2016)









### Nanopore DNA Sequencing









#### What's unique in nanopore sequencing



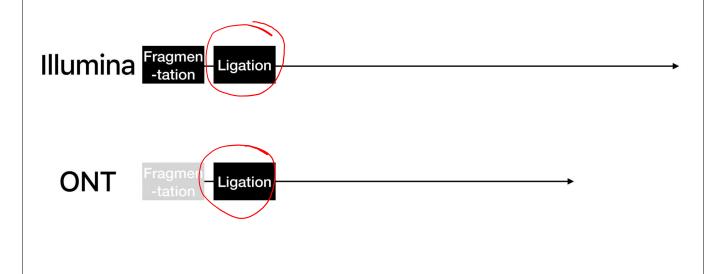
# ONT vs Illumina – Workflow



ONT Fragmen -tation

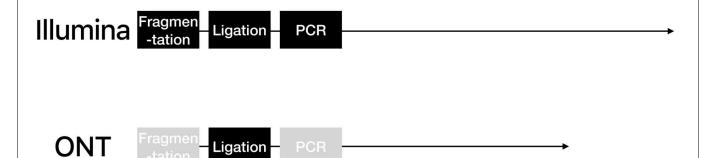


#### **ONT vs Illumina - Workflow**



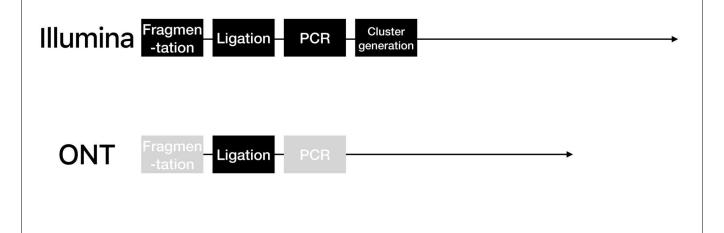


#### **ONT vs Illumina - Workflow**



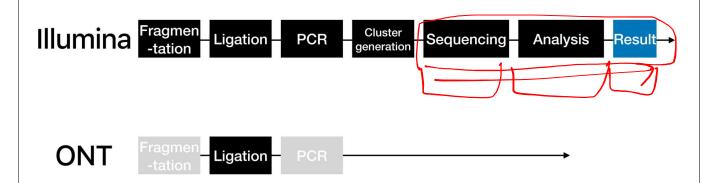


#### **ONT vs Illumina - Workflow**



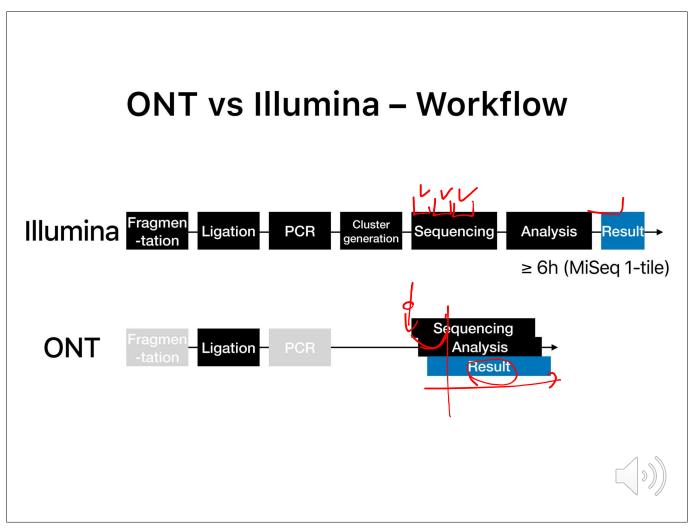


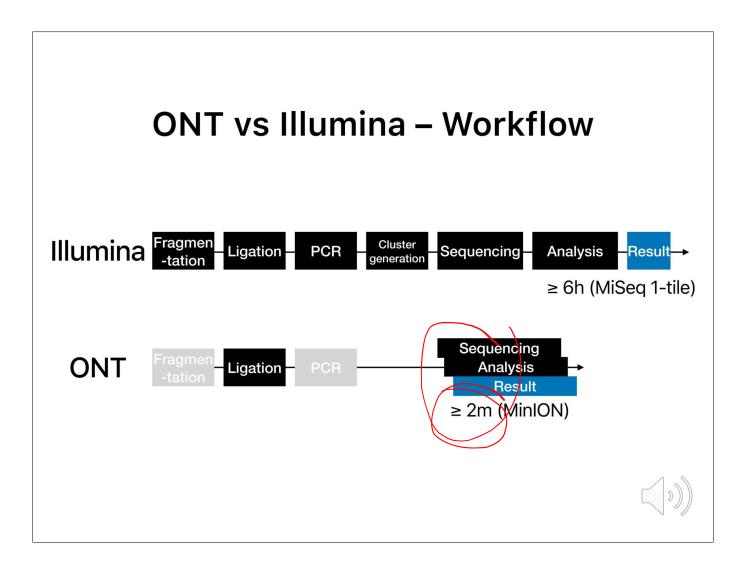
#### **ONT vs Illumina - Workflow**

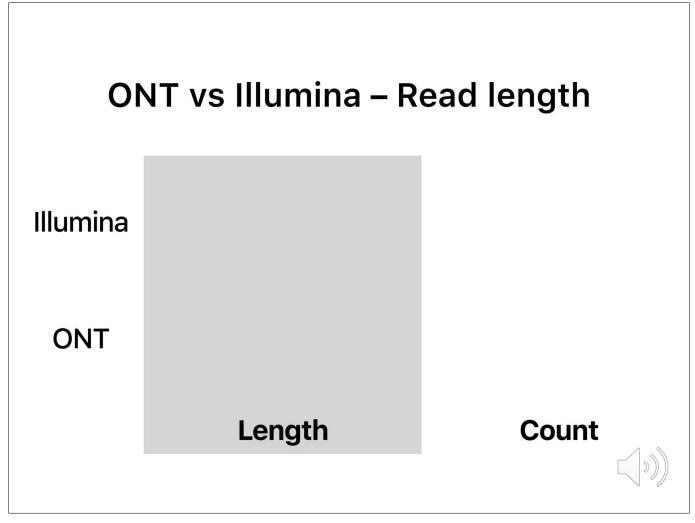




# ONT vs Illumina – Workflow | Cluster generation | Sequencing | Analysis | Result | | 6h (MiSeq 1-tile) | | ONT | Fragmen | Ligation | PCR | PCR | | ONT | Sequencing | Analysis | Result | | 6h (MiSeq 1-tile) | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | Result | | ONT | Sequencing | Analysis | | ONT | ONT | Sequencing | Analysis | | ONT | ONT | Sequencing | Analysis | | ONT | ONT | ONT | | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT | ONT | ONT | ONT | | ONT

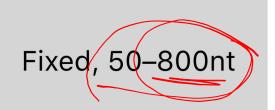






#### ONT vs Illumina - Read length

Illumina



ONT

Length

Count

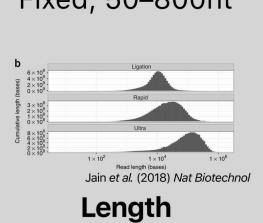


#### **ONT vs Illumina - Read length**

Illumina

Fixed, 50-800nt

ONT



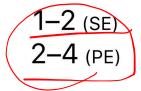
Count



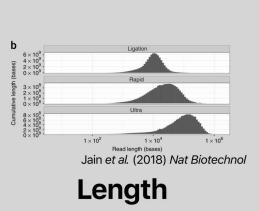
# **ONT vs Illumina - Read length**

Illumina

Fixed, 50-800nt



ONT



Count



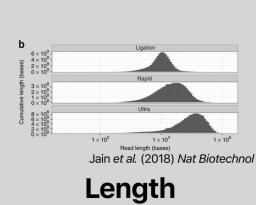
# **ONT vs Illumina - Read length**

Illumina

Fixed, 50-800nt

1-2 (SE) 2-4 (PE)

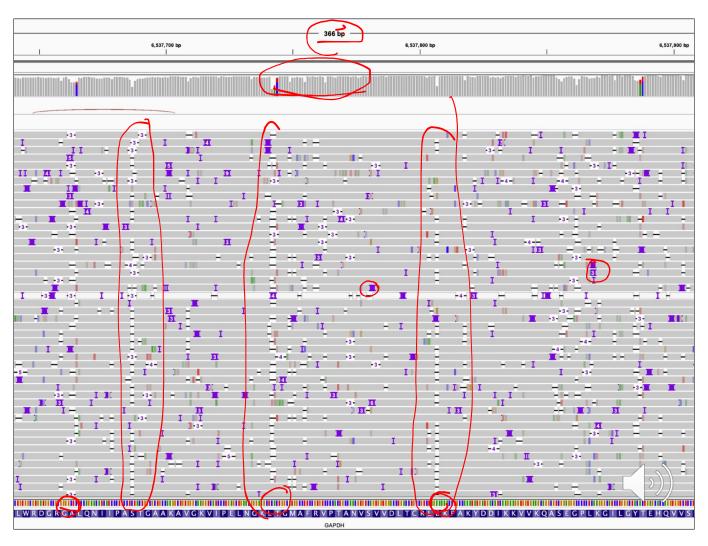
ONT

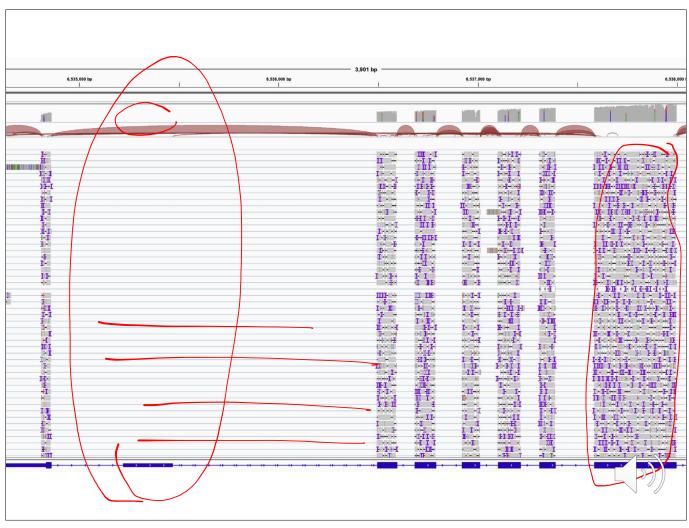


1 (1D) 1-2 (1D<sup>2</sup>) 1-20 (LCS+1D)

Count

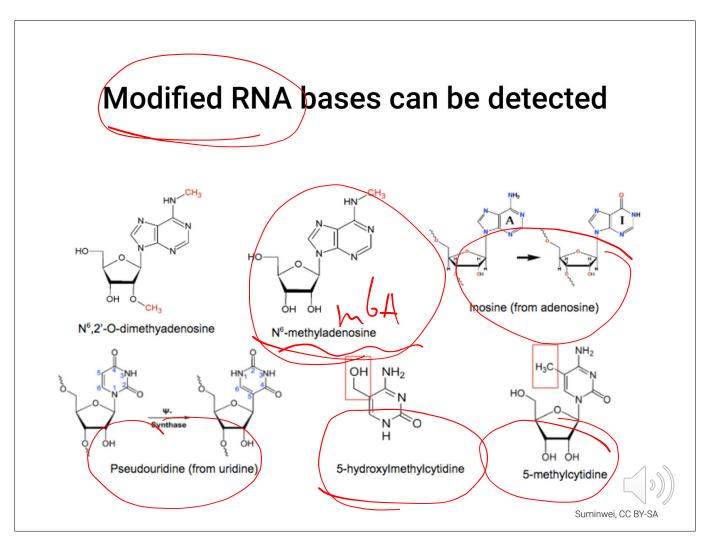


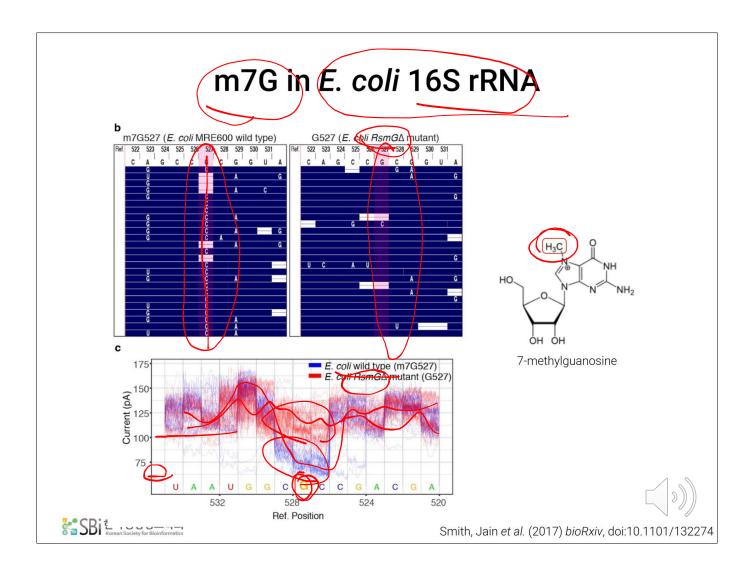


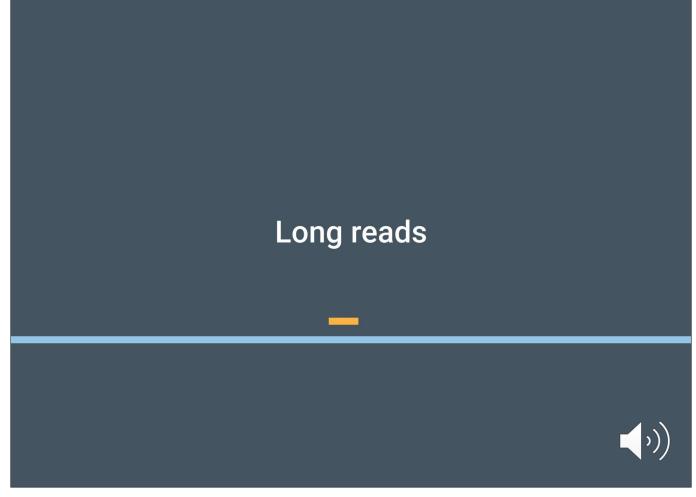


# RNA sequencing with nanopore











"Detects" the splicing isoforms



# "Detects" the splicing isoforms Spans the full length



"Detects" the splicing isoforms

Spans the full length

Read count ~ molecule count



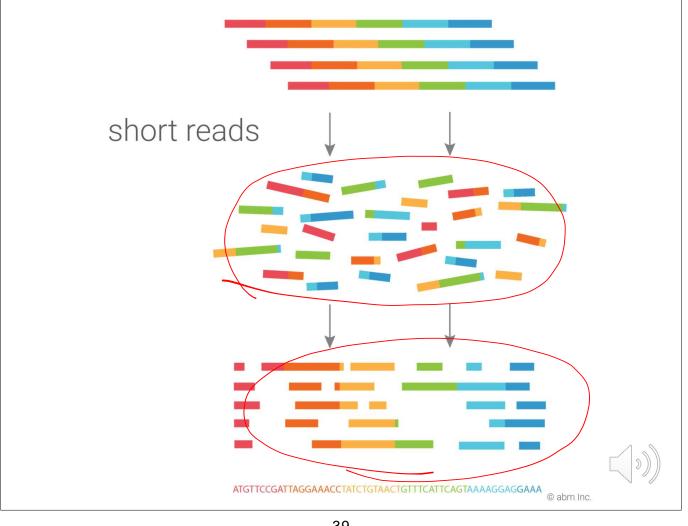
"Detects" the splicing isoforms

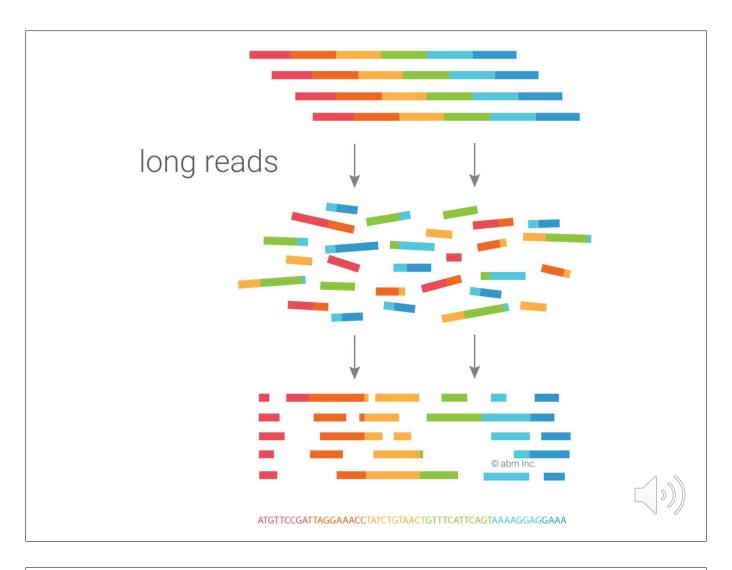
Spans the full length

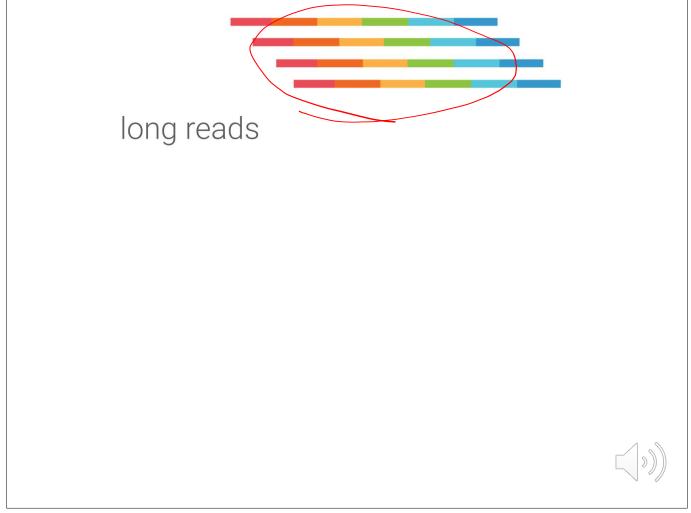
Read count ~ molecule count

Can read long range associations

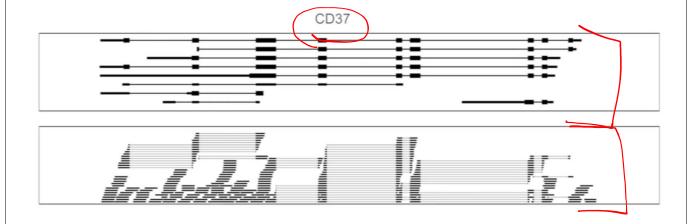








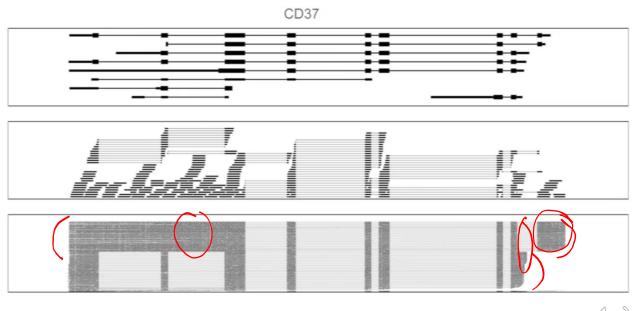
# Long reads enable "reading" of isoforms





Christopher Vollmers (2017) ONT Knowledge Exchange

# Long reads enable "reading" of isoforms



Christopher Vollmers (2017) ONT Knowledge Exchange

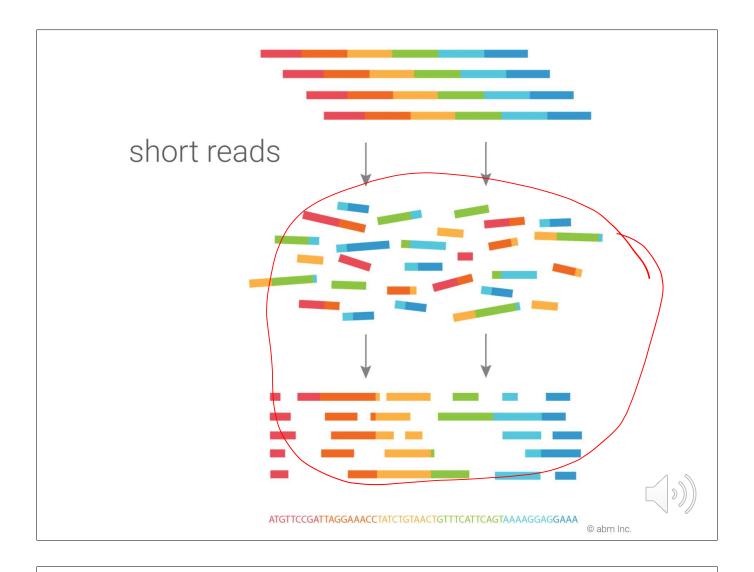


Table 11: Normalization methods for the comparison of gene read counts within the same sample.

### Name

FPKM

TPM

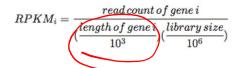
(fragments per

kilobase...

RPKM (reads per kilobase of exons per million mapped reads)

### Details

- 1. For each gene, count the number of reads mapping 2. Divide that count by: the length of the gene in
- base pairs divided by 1,000 multiplied by the total number of mapped reads divided by  $10^6$ .



- 1. Same as RPKM, but for paired-end reads:
  - 2. The number of fragments (defined by two reads each) is used.

Instead of normalizing to the total library size, TPM represents the abundance of an individual gene i in relation to the abundances of the other transcripts (e.g.,

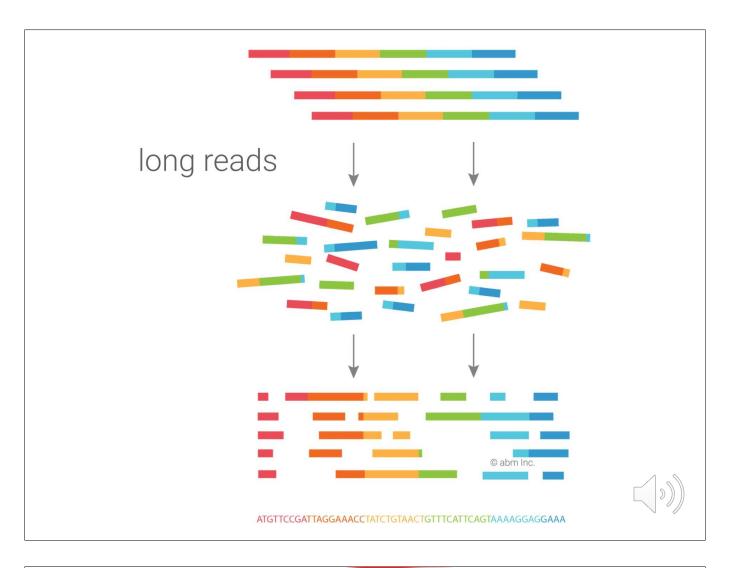
- 1. For each gene, count the number of reads mapping to it and divide by its length in base pairs (= counts per base).
- 2. Multiply that value by 1 divided by the sum of all counts per base of every gene.
- Multiply that number by 10<sup>6</sup>.

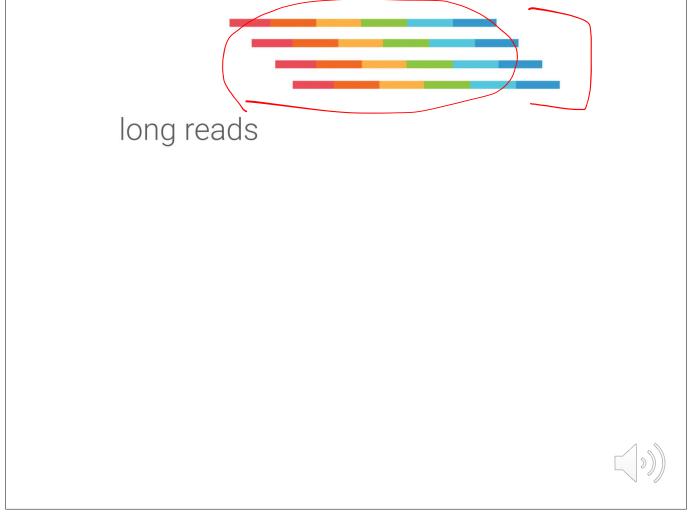
 $TPM_i = \frac{X_i}{l_i} *$ 

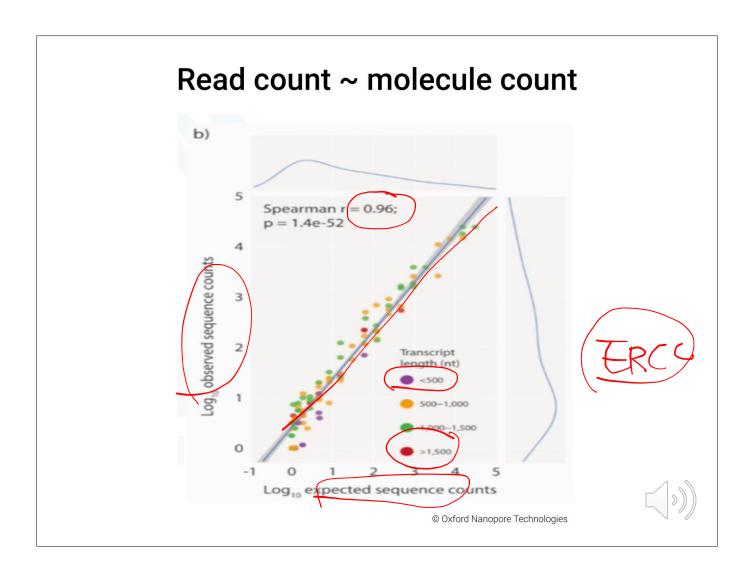
### Comment

- · introduces a bias in the per-gene variances, in particular for lowly expressed genes (Oshlack and Wakefield, 2009)
- implemented in edgeR's rpkm() function
- implemented in DESeq2's fpkm() function
- details in Wagner et al. (2012)

Dünbar et al. (2015)





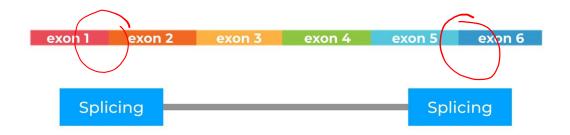






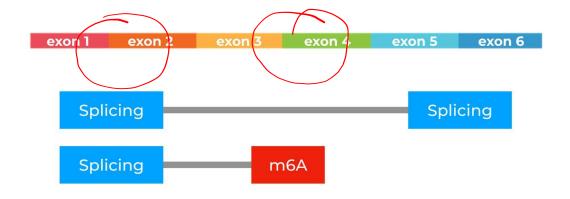


# Capturing long-range associations

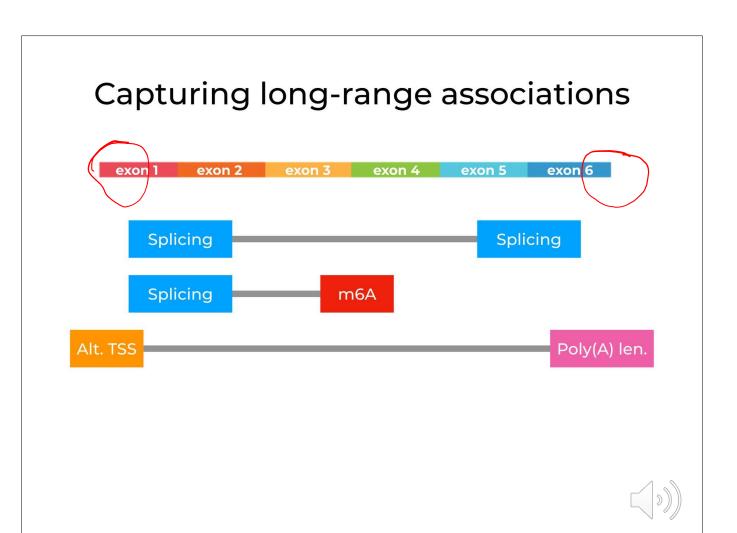


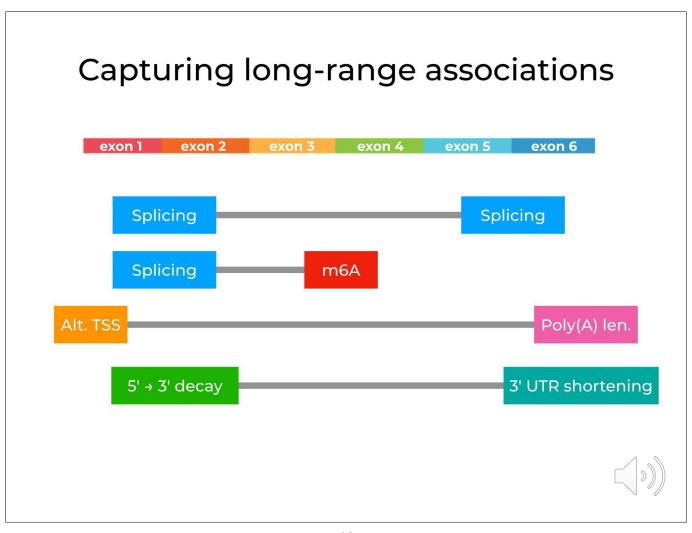


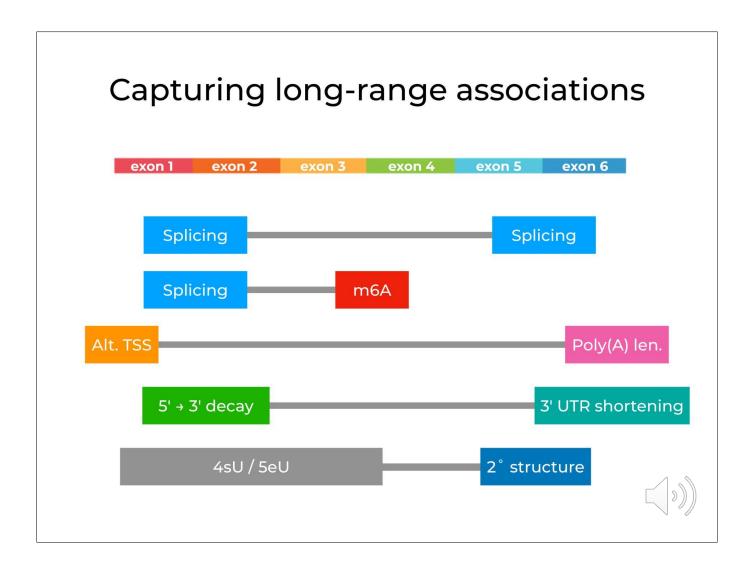
# Capturing long-range associations

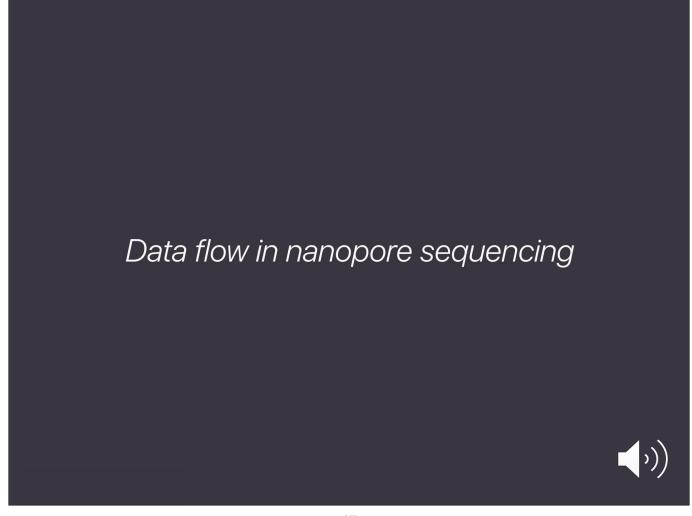


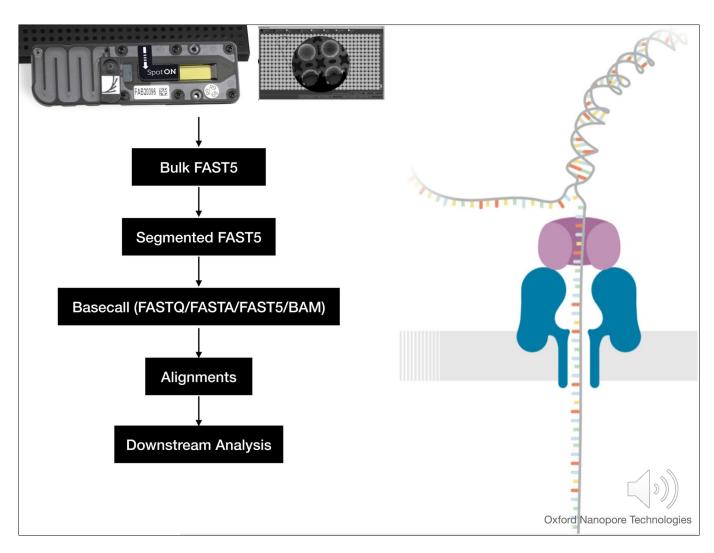


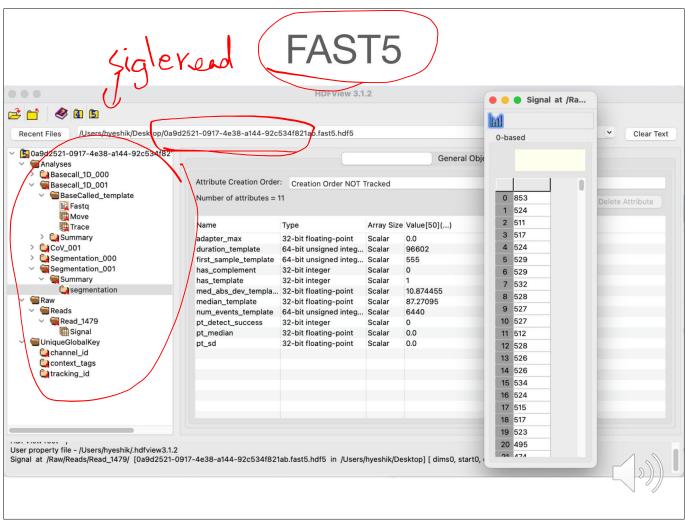


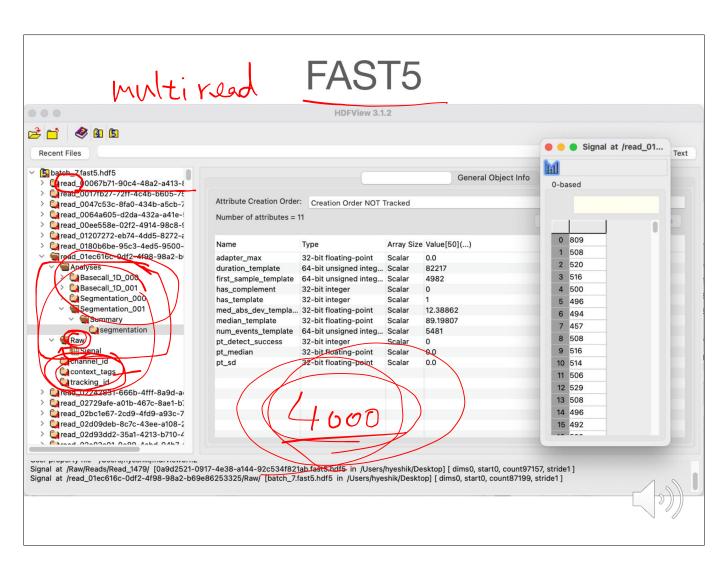


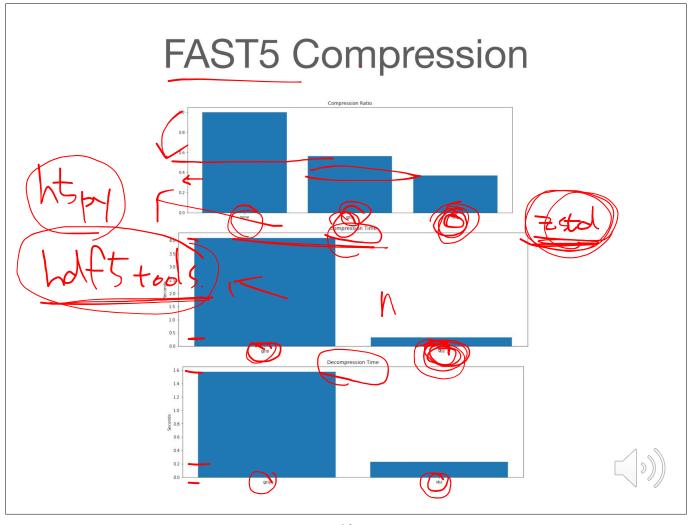








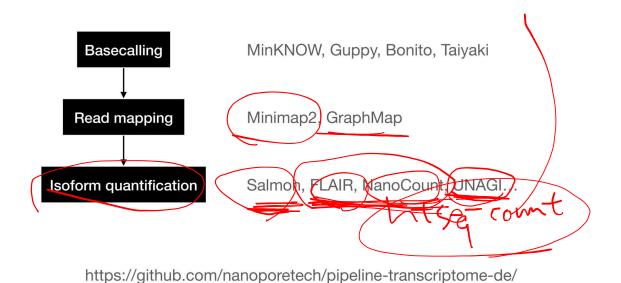


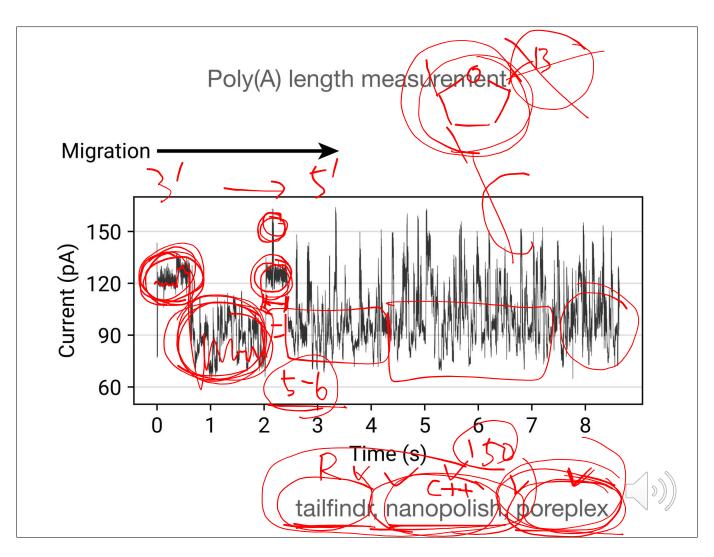


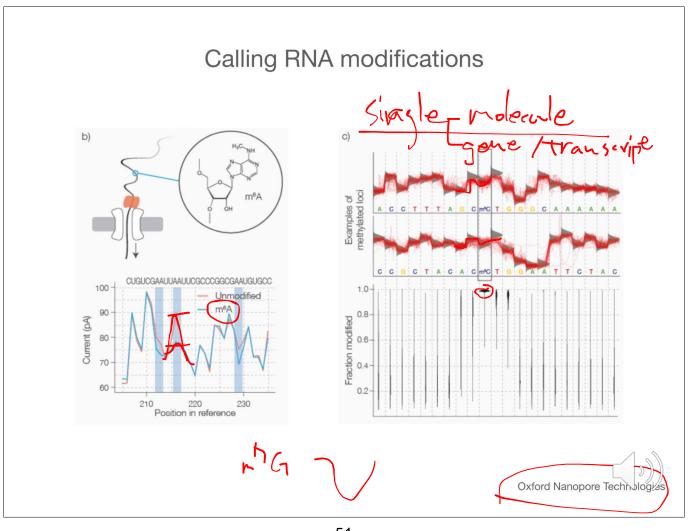
# Practical workflows and tools for expression analysis using nanopore sequencing



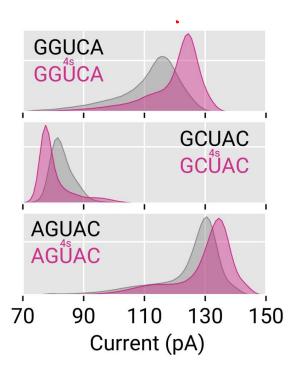
### Expression profiling using nanopore sequencing



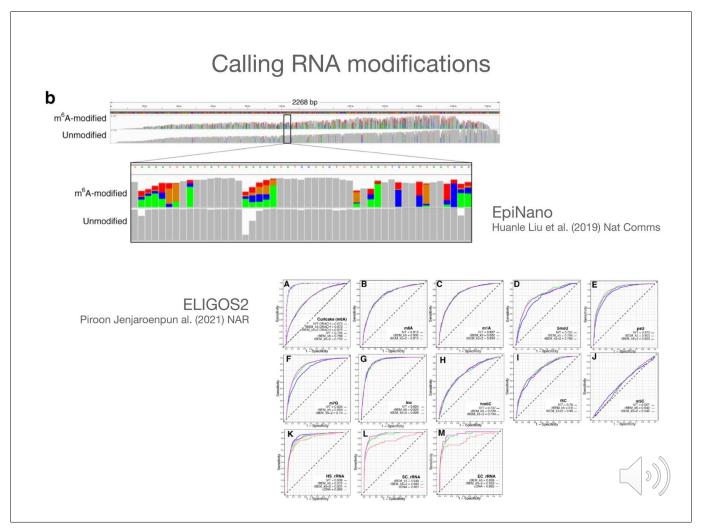


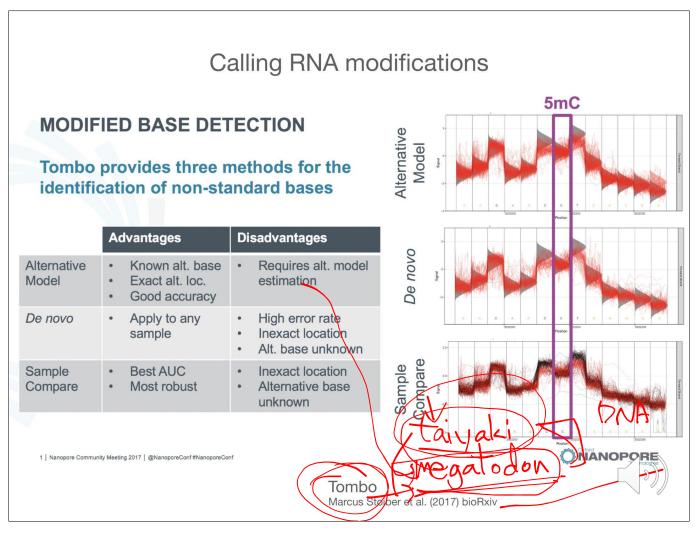


### Calling RNA modifications

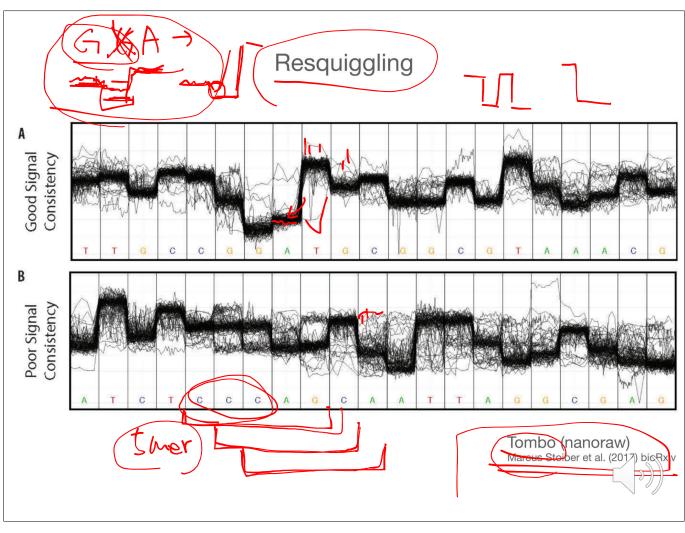


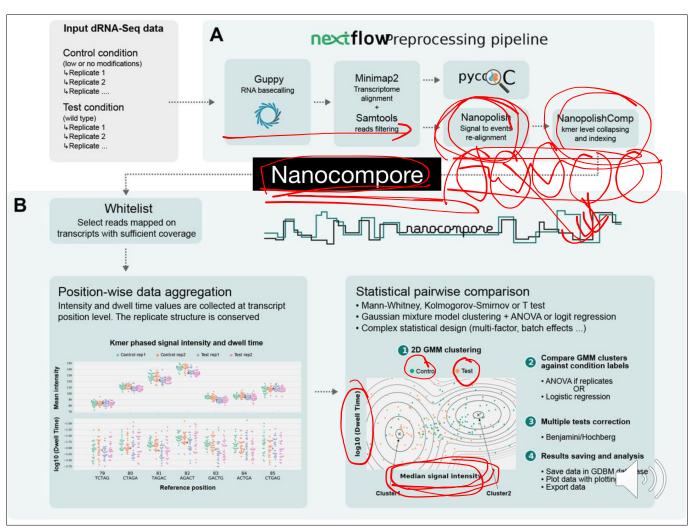




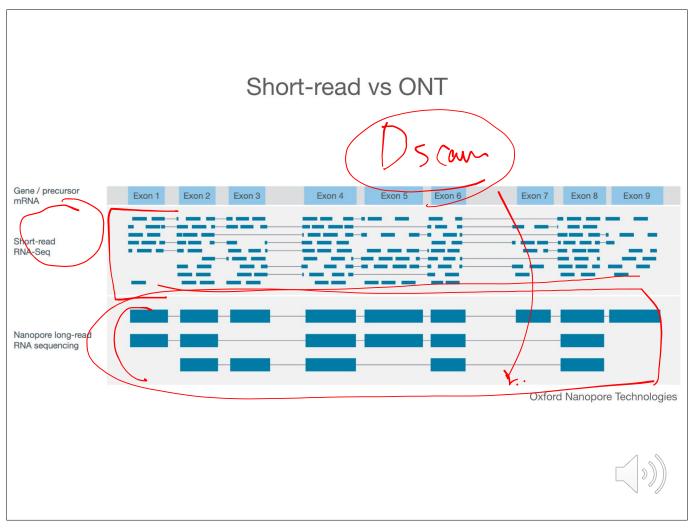


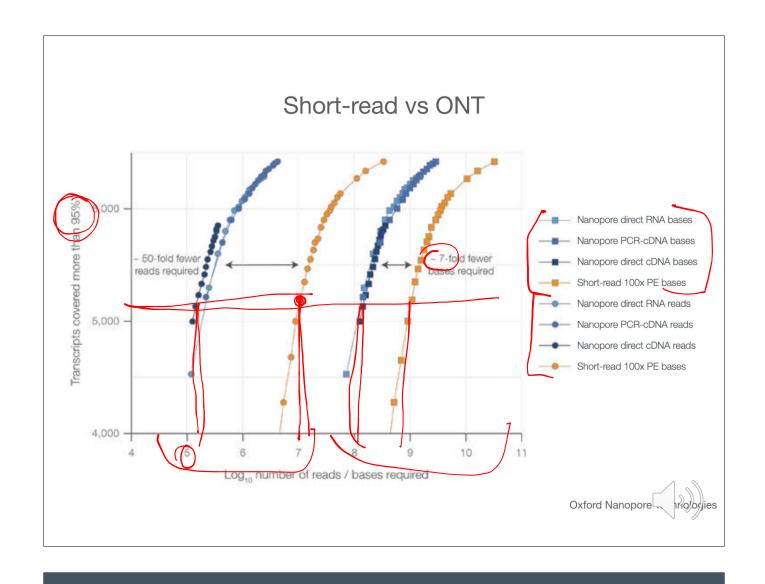












# **DO NOT USE ONT direct RNA sequencing**



# DO NOT USE ONT direct RNA sequencing (yet) for:



# **DO NOT USE ONT direct RNA sequencing**

(yet) for:

• small RNA (miRNA, snoRNA, etc.)



# DO NOT USE ONT direct RNA sequencing (yet) for:

- small RNA (miRNA, snoRNA, etc.)
- Poly(A) RNAs with variable 3' ends



# DO NOT USE ONT direct RNA sequencing (yet) for:

- small RNA (miRNA, snoRNA, etc.)
- Poly(A) RNAs with variable 3' ends
- <100 ng input RNA</p>



# DO NOT USE ONT direct RNA sequencing (yet) for:

- small RNA (miRNA, snoRNA, etc.)
- Poly(A) RNAs with variable 3' ends
- <100 ng input RNA</p>
- Study objective requires very high baselevel quality





Quantification of poly(A)+ RNAs



- Quantification of poly(A)+ RNAs
- Isoform-level analyses

- Quantification of poly(A)+ RNAs
- Isoform-level analyses
- When without a high quality reference



- Quantification of poly(A)+ RNAs
- Isoform-level analyses
- When without a high quality reference
- Detecting modified bases



- Quantification of poly(A)+ RNAs
- Isoform-level analyses
- When without a high quality reference
- Detecting modified bases
- Discovery of long-range associations



- Quantification of poly(A)+ RNAs
- · Isoform-level analyses
- When without a high quality reference
- Detecting modified bases
- Discovery of long-range associations
- Measuring poly(A) lengths



- Quantification of poly(A)+ RNAs
- Isoform-level analyses
- When without a high quality reference
- Detecting modified bases
- Discovery of long-range associations
- Measuring poly(A) lengths
- · When a fast feedback is desirable

