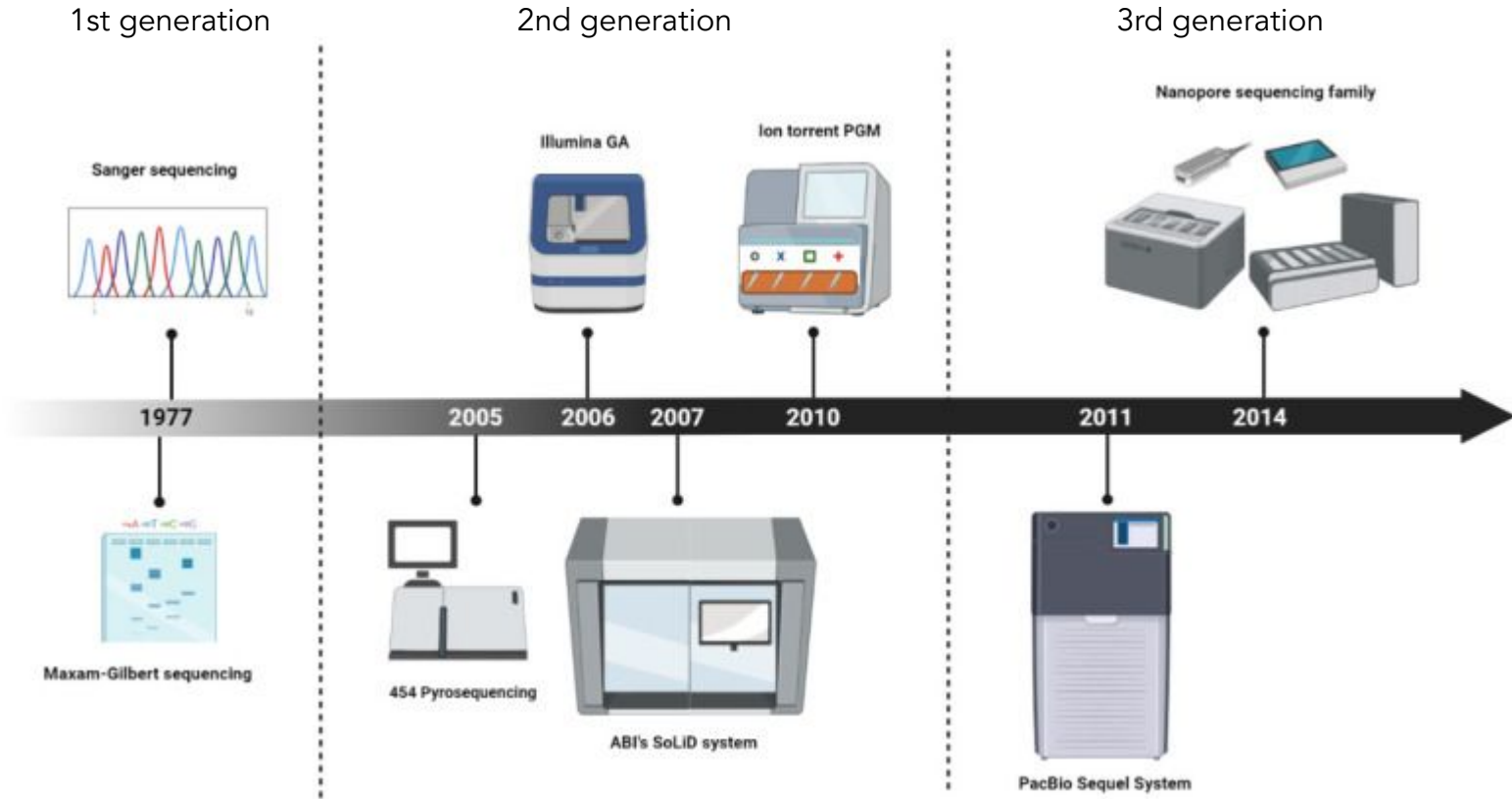


Single-molecule sequencing

Quantitative Biology 2022
10/28/22

The history of sequencing



Single-molecule sequencing (SMS)



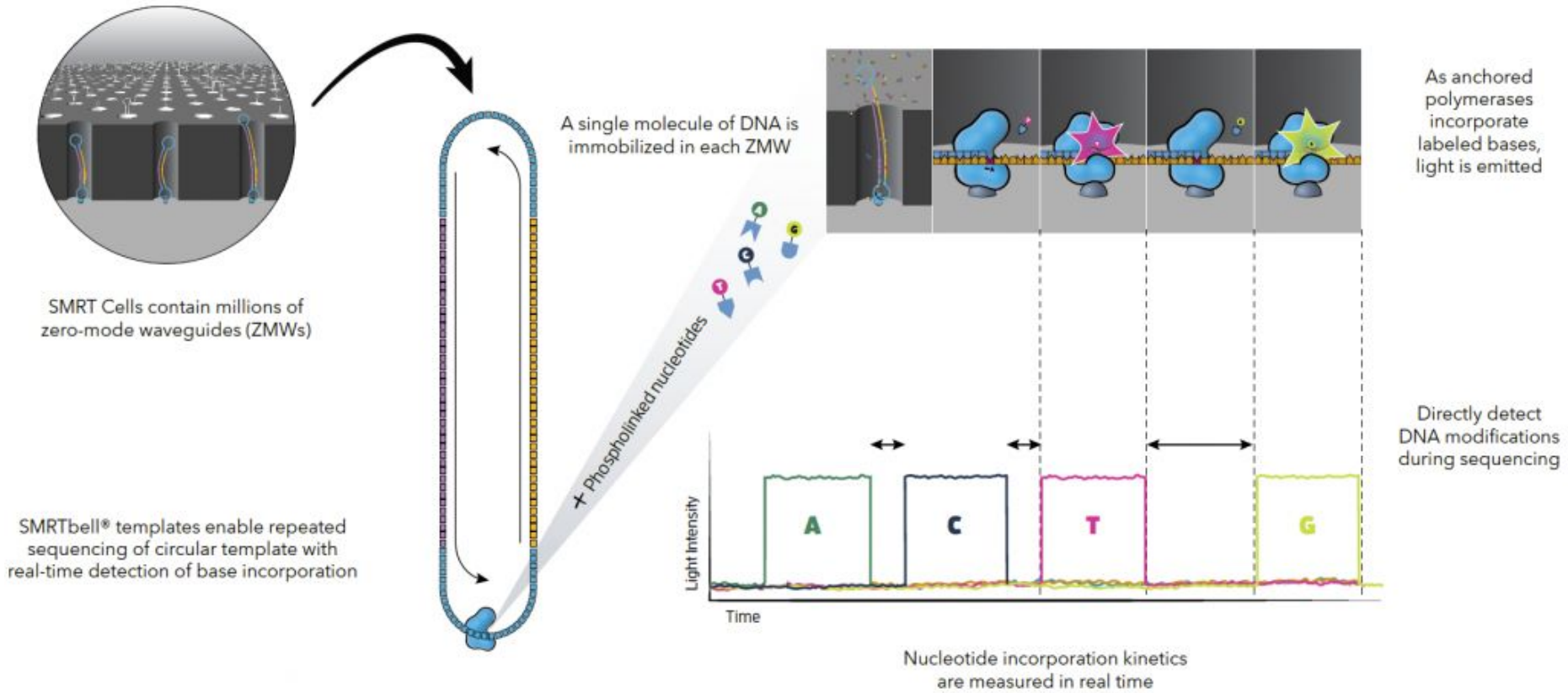
Advantages

- No PCR
- Small input sample size
- Real-time sequencing
- More uniform genome coverage
- Longer read lengths (Kb to Mb)
- Faster data production

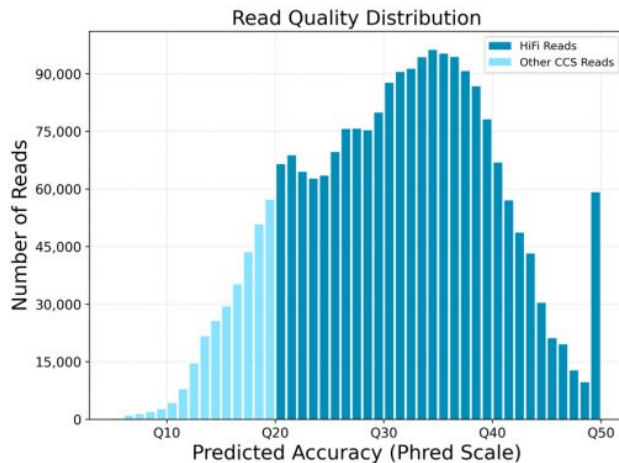
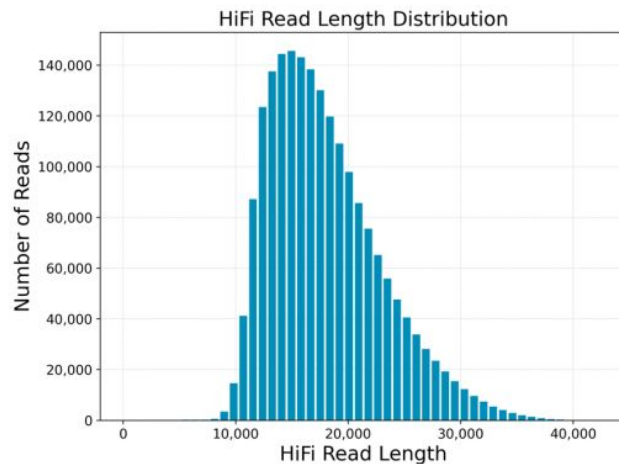
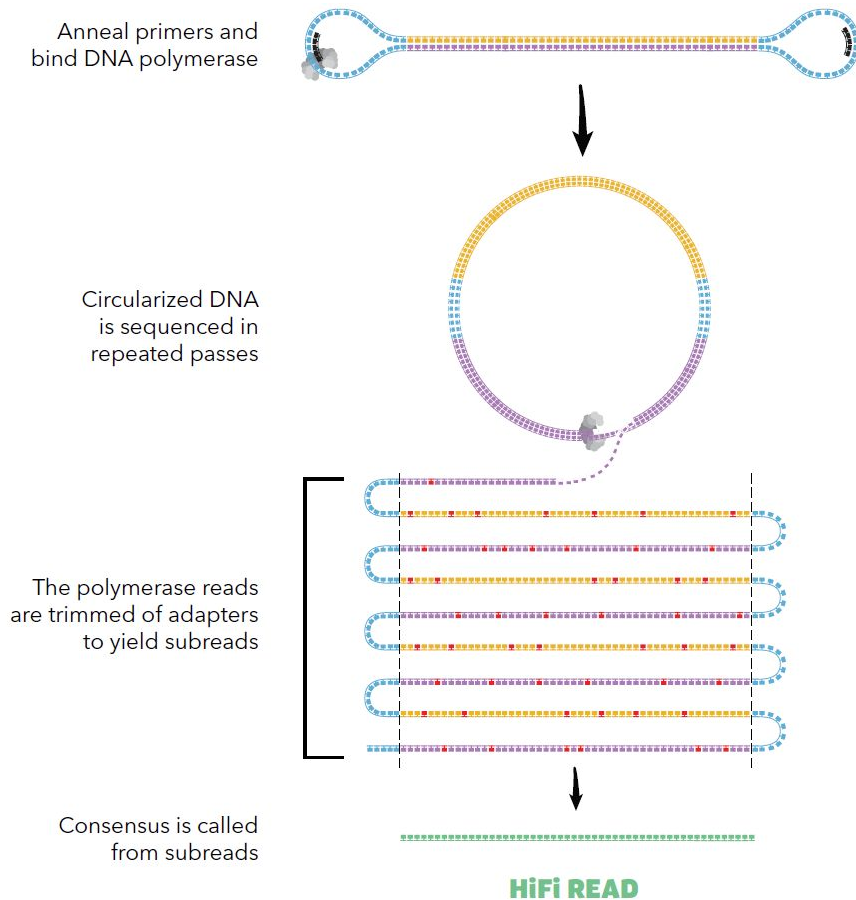
Disadvantages

- Higher error rates
- Fresh samples and careful handling needed to preserve ultralong reads
- Lack of database/analysis tools

PacBio - sequencing



PacBio - basecalling



PacBio

Advantages

- No PCR
- Small input sample size
- Real-time sequencing
- More uniform genome coverage
- Longer read lengths (up to 300Kb, ~15Kb average)
- Faster data production

Disadvantages

- Higher error rate (~90% accurate, HiFi ~99.9% accurate)
- Requires fresh samples
- Immature database/analysis tools

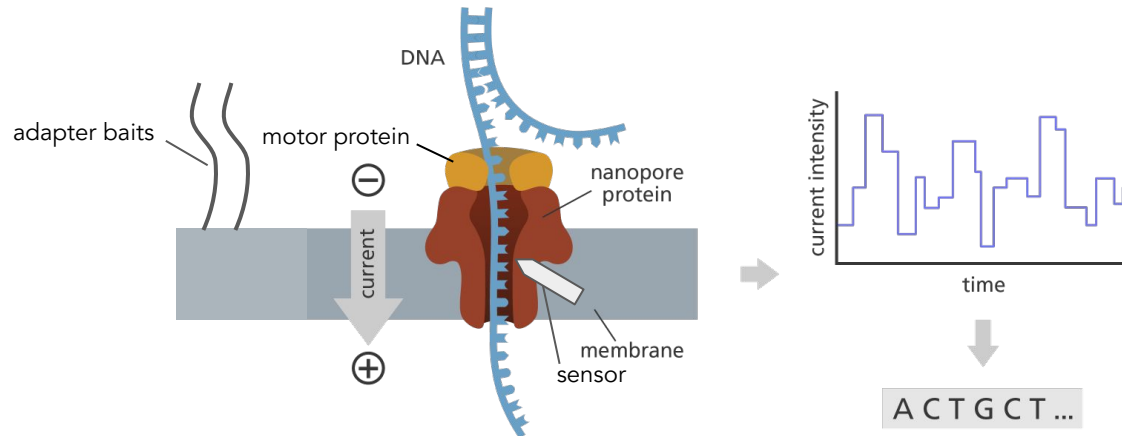
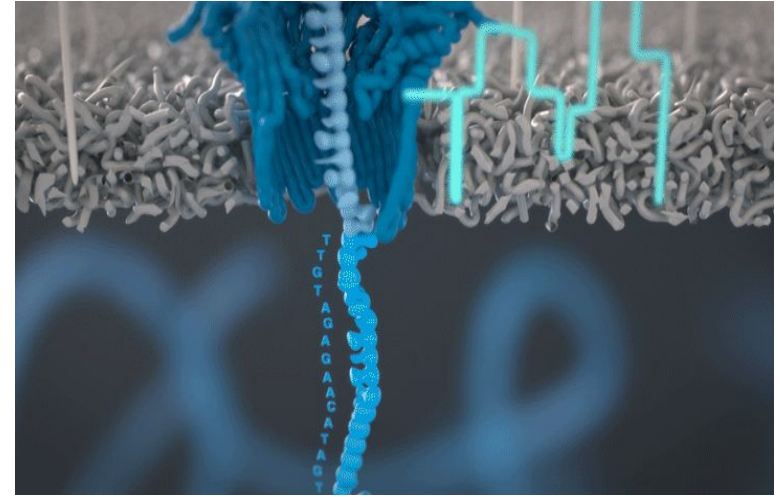
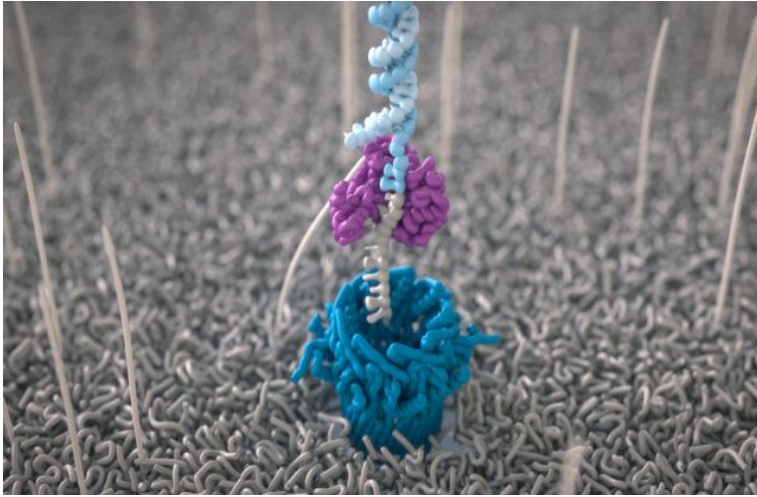
PacBio

A single run on a PacBio Sequel IIe (~\$500K/machine)

- Yields ~4 million reads
- Takes about 30 hours
- Sequences fragments with an average size of 15Kb
- Generates ~60Gbp of sequence
- Costs ~\$2-3K

For comparison, an Illumina run for the same amount of data would cost ~\$1-1.5K

Oxford Nanopore Technologies (ONT)



Nanopore sequencing

Advantages

- No PCR
- Portable
- Low equipment cost
- Can directly sequence RNA
- Small input sample size (~order magnitude less than PacBio)
- Real-time sequencing
- More uniform genome coverage
- Longer read lengths (up to 4Mb, ~50-100Kb average)
- Faster data production
- Reusable flowcells

Disadvantages

- Higher error rate (90-95% accurate)
- Requires fresh samples
- Immature database/analysis tools

Nanopore Sequencing

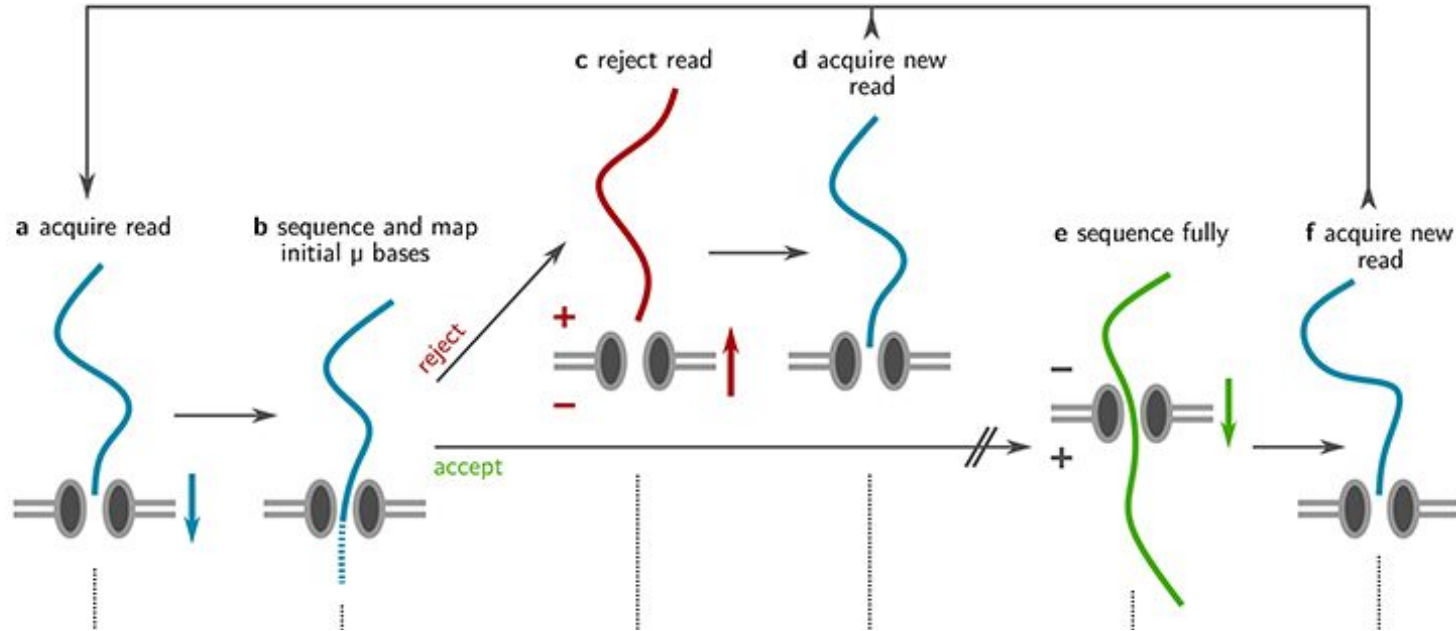
A single run on a Minlon (~\$1K/machine)

- Yields ~200 thousand reads
- Takes about 24-72 hours
- Sequences fragments with an average size of 50-100Kb
- Generates ~10-20Gbp of sequence
- Costs ~\$500

For comparison, an Illumina run for the same amount of data would cost
~\$300

ONT advances

"Read until" Adaptive Read Selection

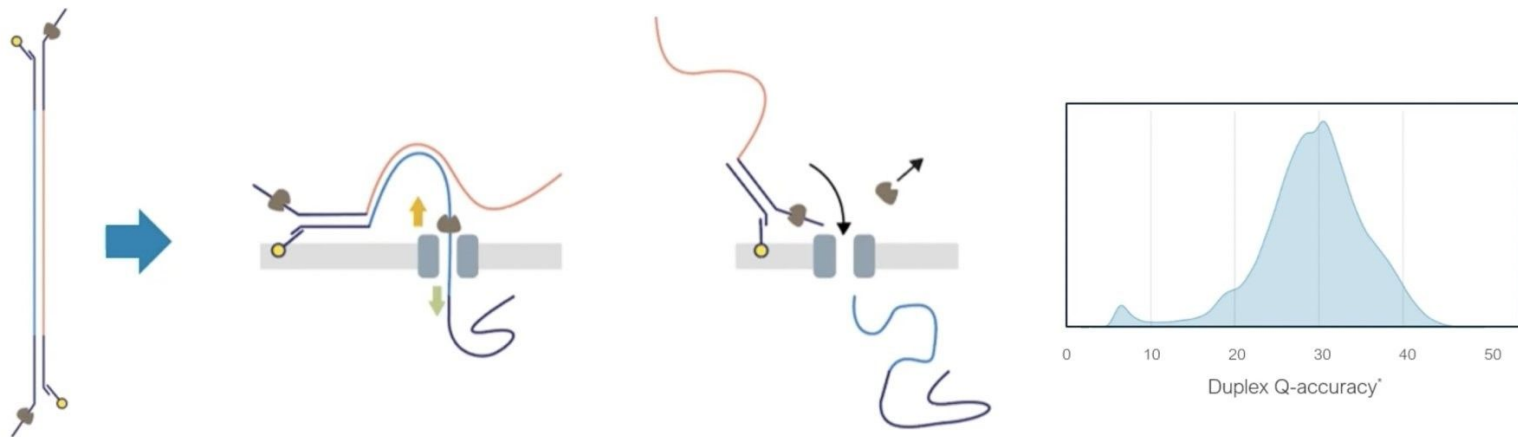


Sequence-based read acceptance/rejection allows enrichment of target templates at the time of sequencing

ONT advances



Duplex Sequencing

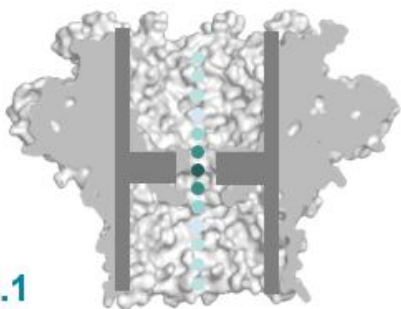


Forward-reverse sequencing allows joint basecalling, vastly improving read accuracy (~99.9% accurate on average)

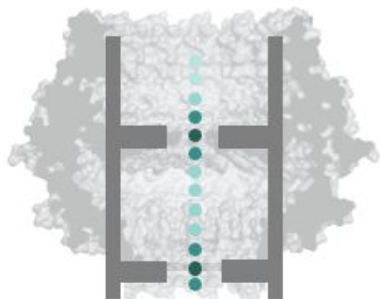
ONT advances

Dual Pore Sensors

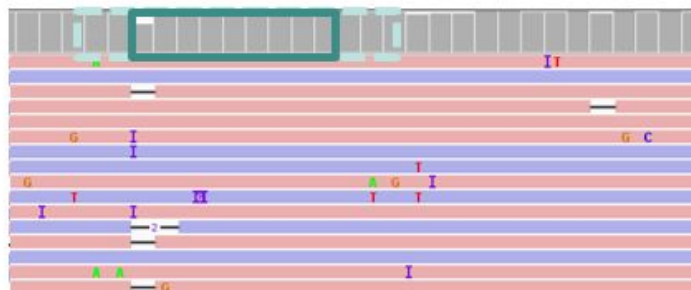
R9.4.1



R10



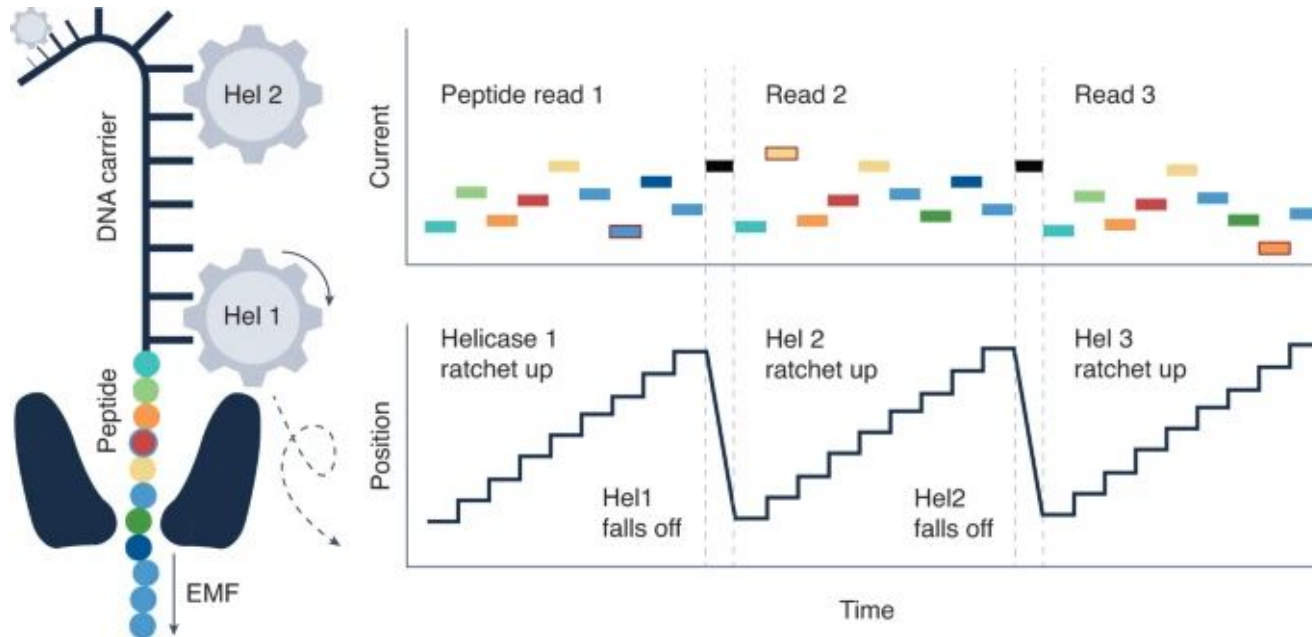
ATCGGAAAAAAAAAATCACGCCACGTCCAAA



Multiple offset sensors improve resolution of homopolymers

ONT advances

Direct peptide sequencing

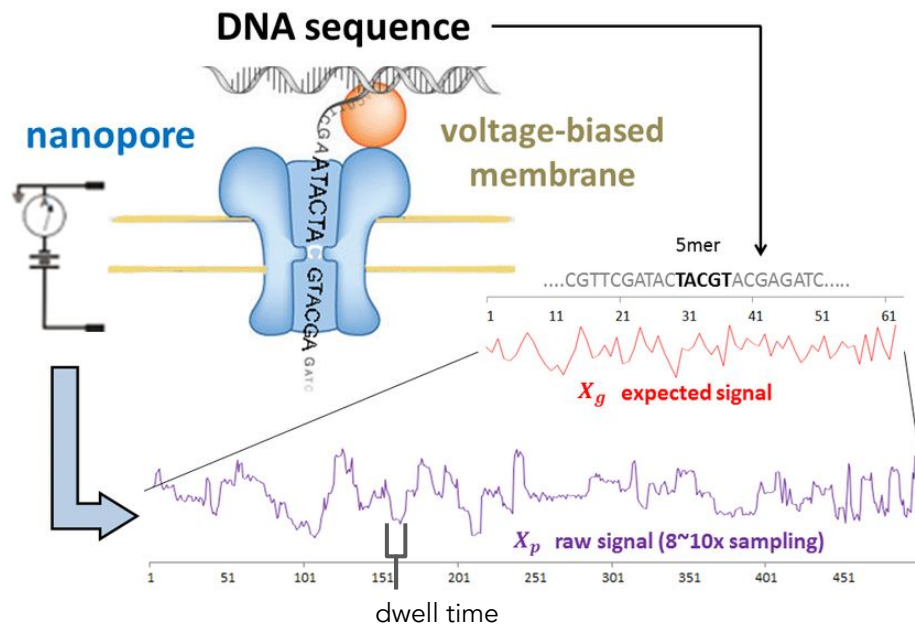


Helicases ratchet peptides through pores repeatedly for a high quality consensus protein sequence

From voltage to base calls

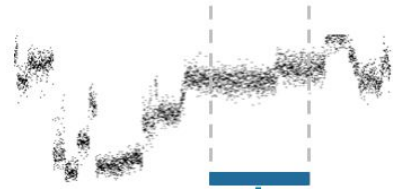


- Voltage is determined by 5 bases currently in pore
- Each 5-mer yields (semi-)unique voltage
- Dwell time in pore is variable



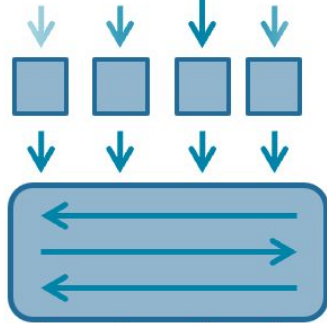
Guppy - a recurrent neural network basecaller

Base calling (RNN, raw)



Parameters learned from training data

Extraction of blocks of features



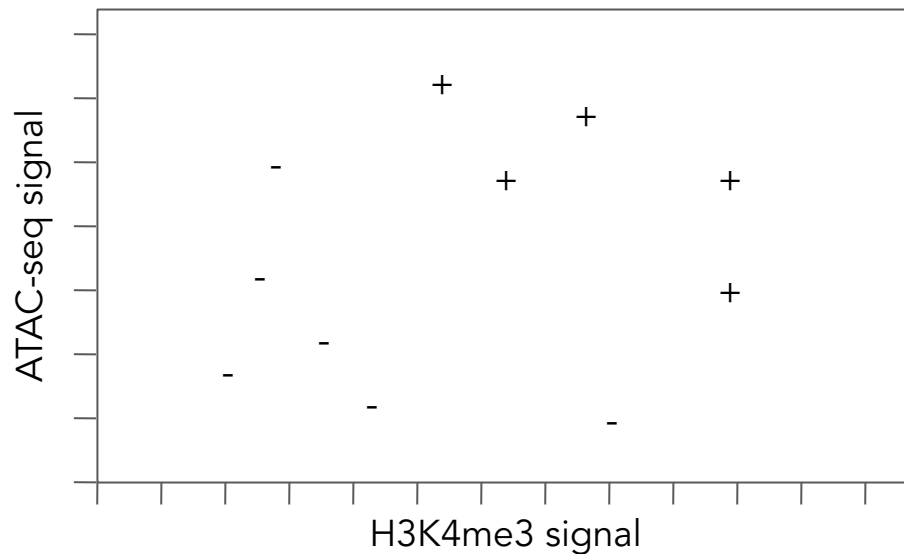
Bidirectional
information
flow

Multi-base prediction

Decode to sequence

What is a neural network

Let's take a step back and look at the simple case of 2 category classification

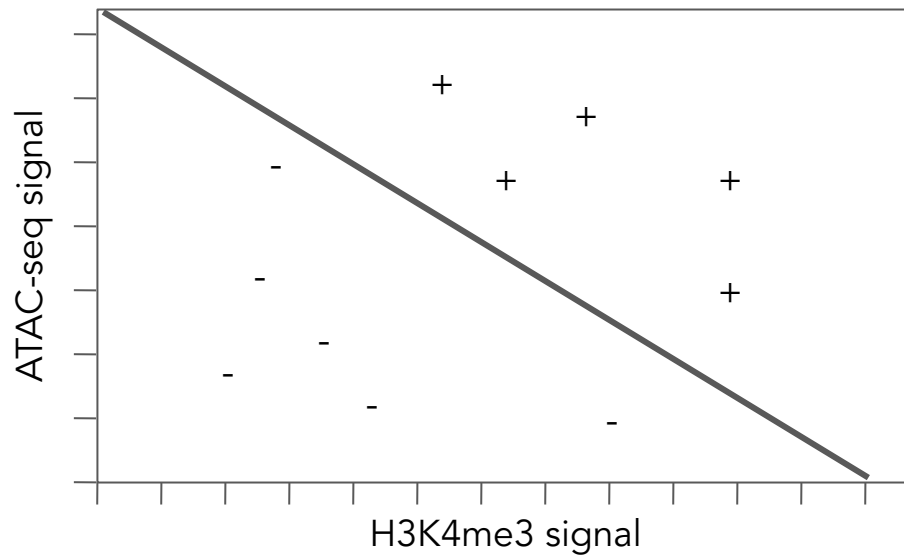


We want to determine if a likely active using H3K4me3 and ATAC-seq signal

We start with known examples

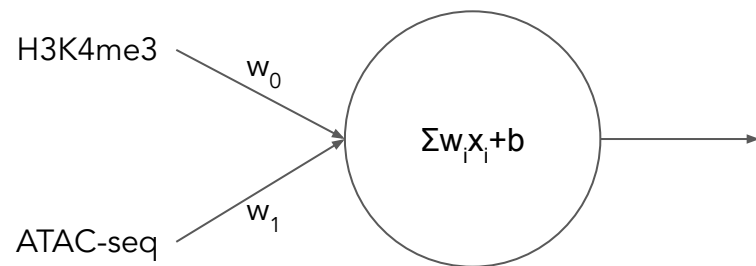
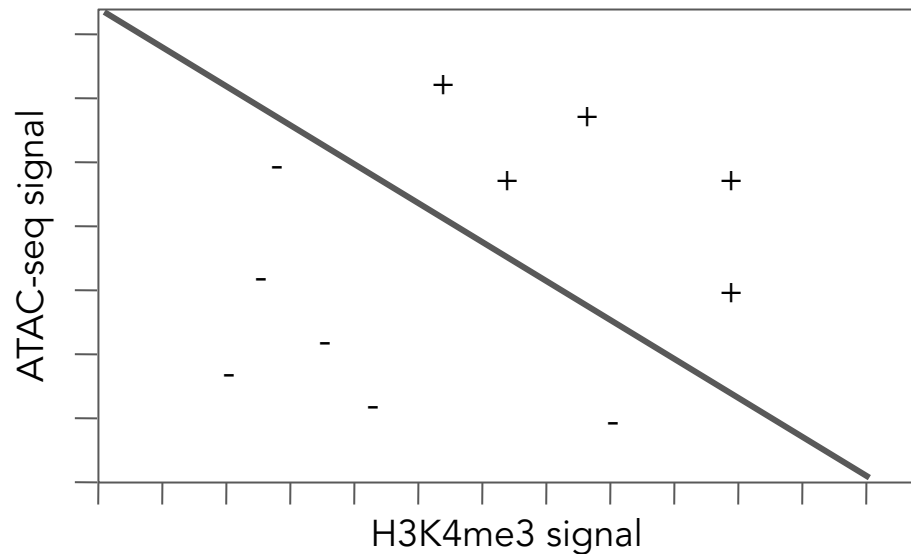
What is a neural network

This is easy to classify with a simple linear relationship



What is a neural network

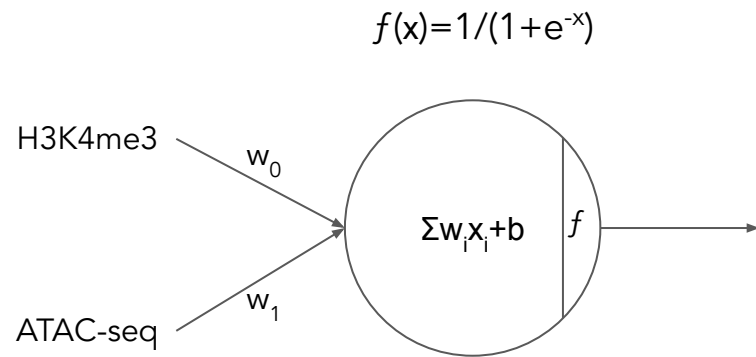
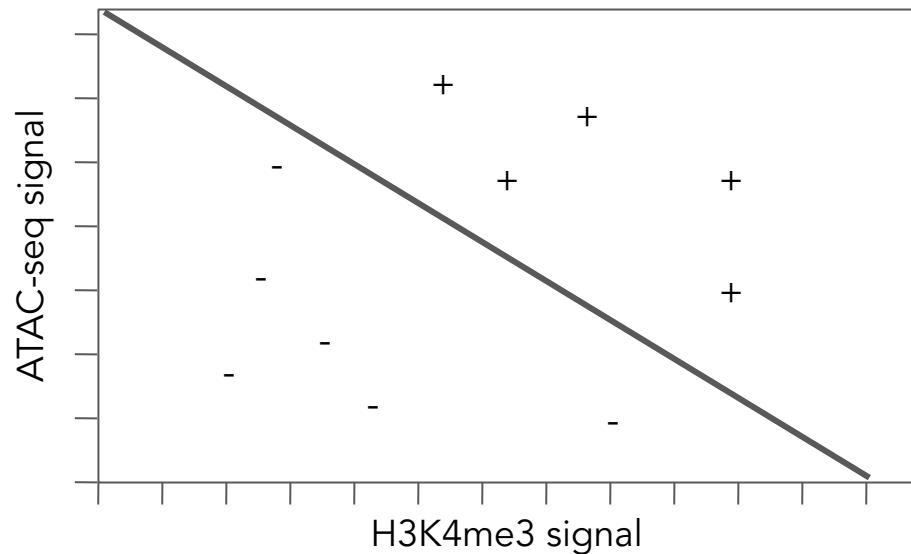
This is easy to classify with a simple linear relationship



We can create a simple linear combination (regression) of the inputs to define the line

What is a neural network

This is easy to classify with a simple linear relationship

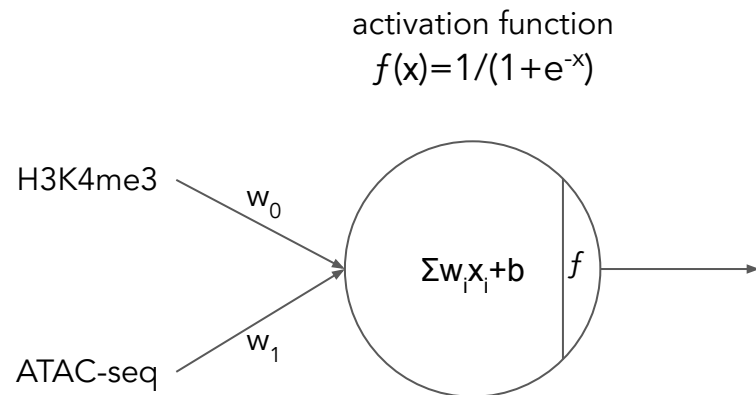
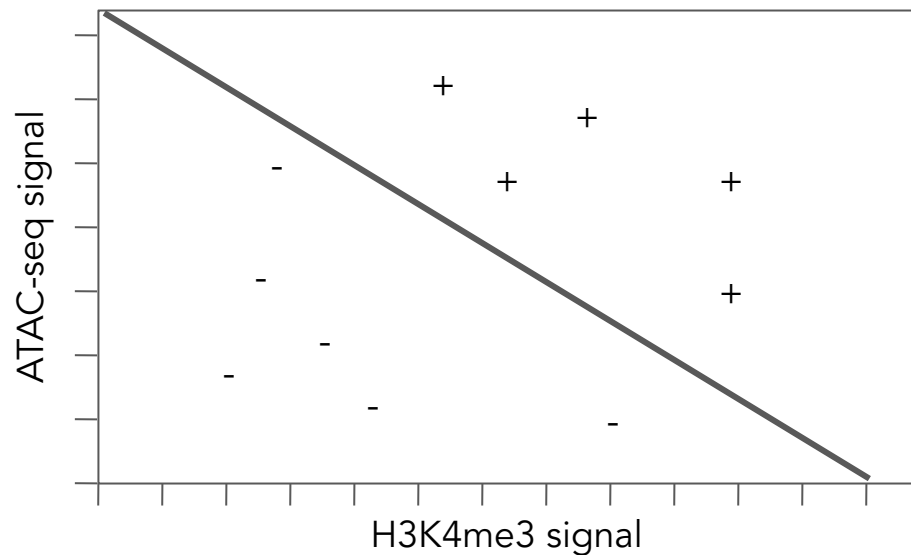


We can create a simple linear combination (regression) of the inputs to define the line

We can add a sigmoid function to change the output to 0-1. This turns it into a logistic regression.

The perceptron

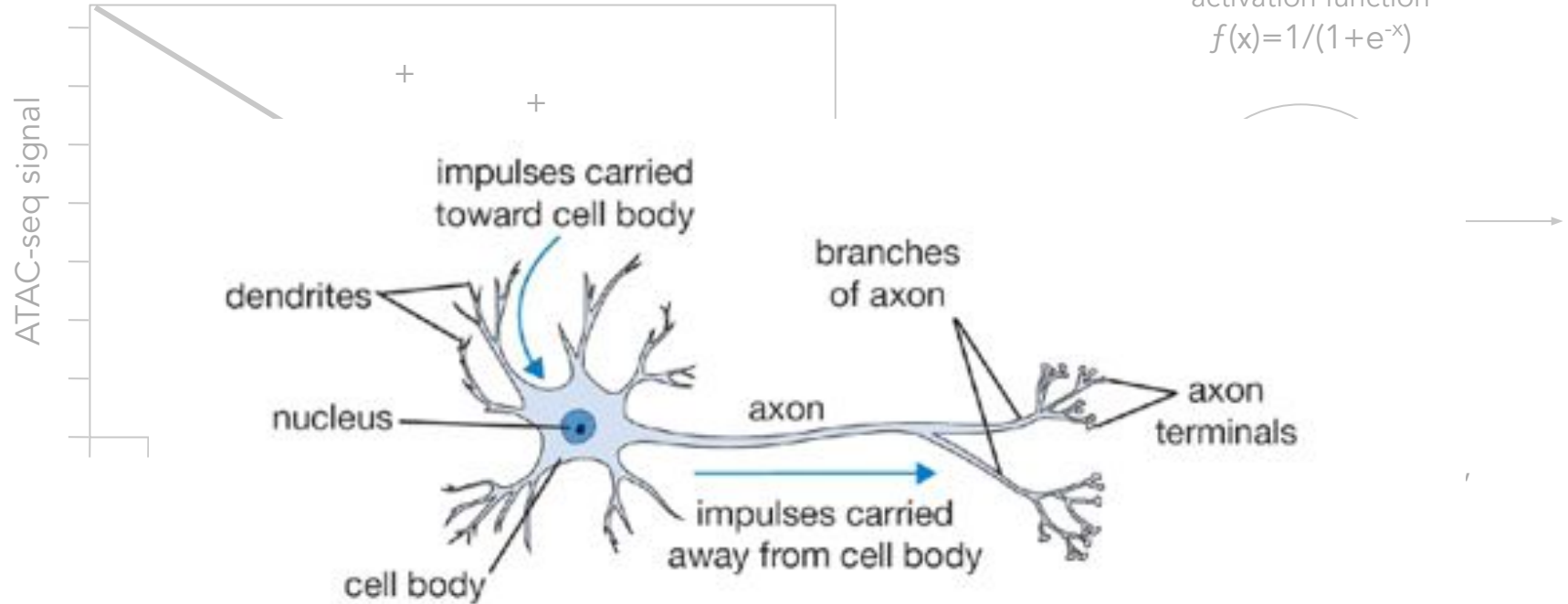
This is easy to classify with a simple linear relationship



We have just defined a “perceptron”

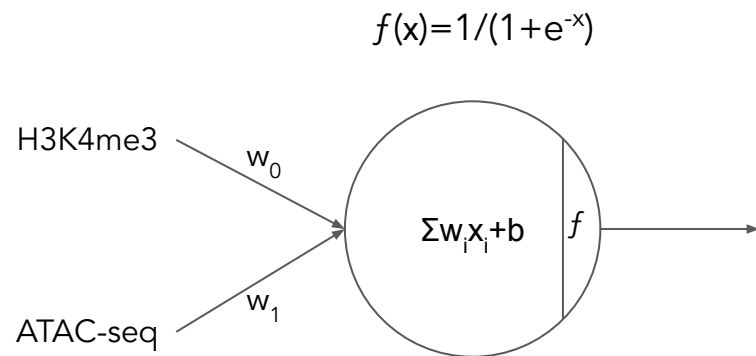
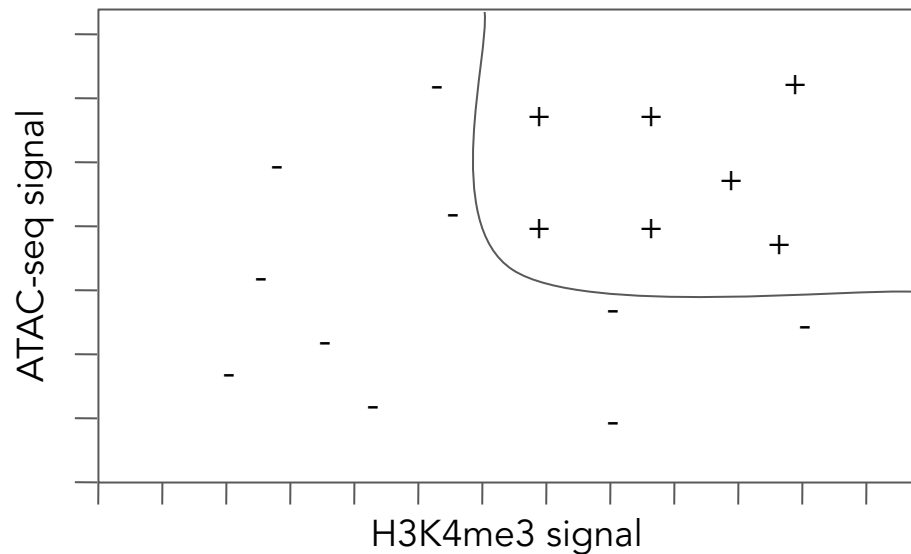
The “neural” in neural network

This is easy to classify with a simple linear relationship



When a perceptron isn't enough

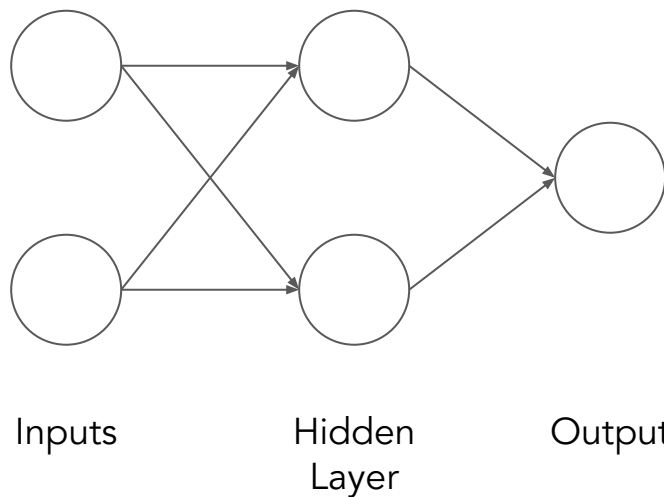
What if it's not a linear relationship?



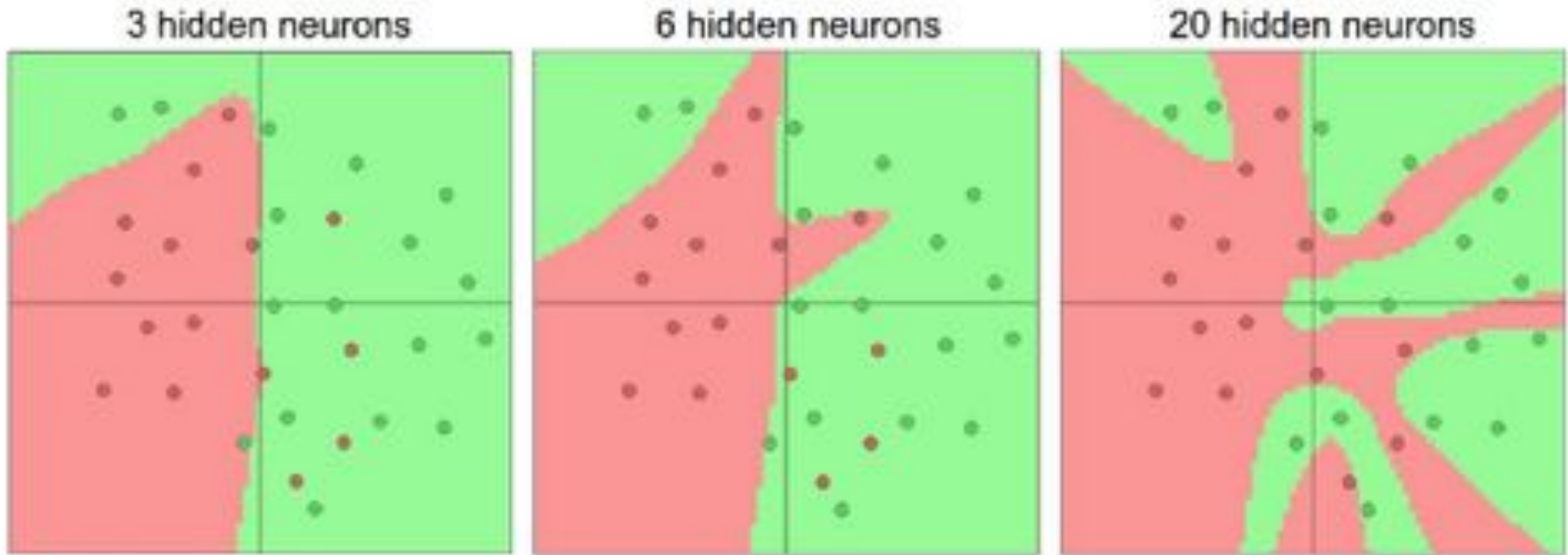
A simple neural network

By combining multiple perceptrons, we can create non-linear functions.

This is a 2-layer network (inputs are not counted)



A universal approximator

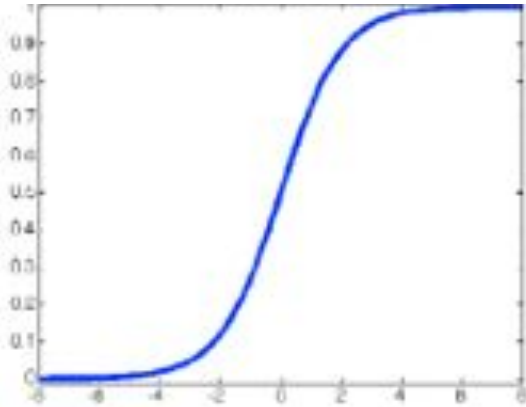


Given enough hidden nodes (neurons), any nonlinear function can be approximated

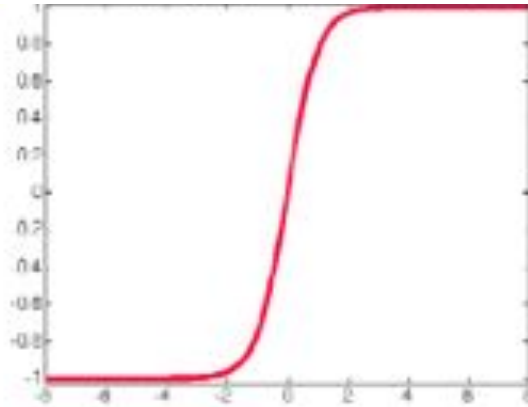
Activation functions



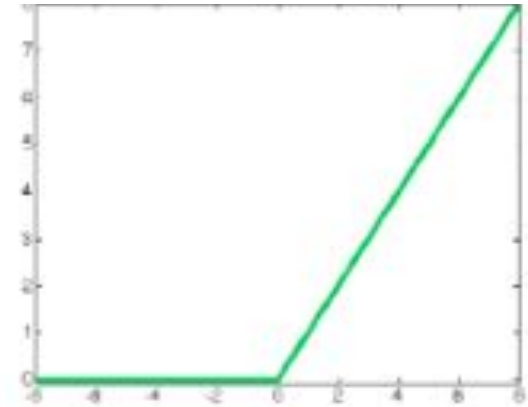
Sigmoid: $f(x) = 1/(1+e^{-x})$



Tanh: $f(x) = (e^x - e^{-x}) / (e^x + e^{-x})$

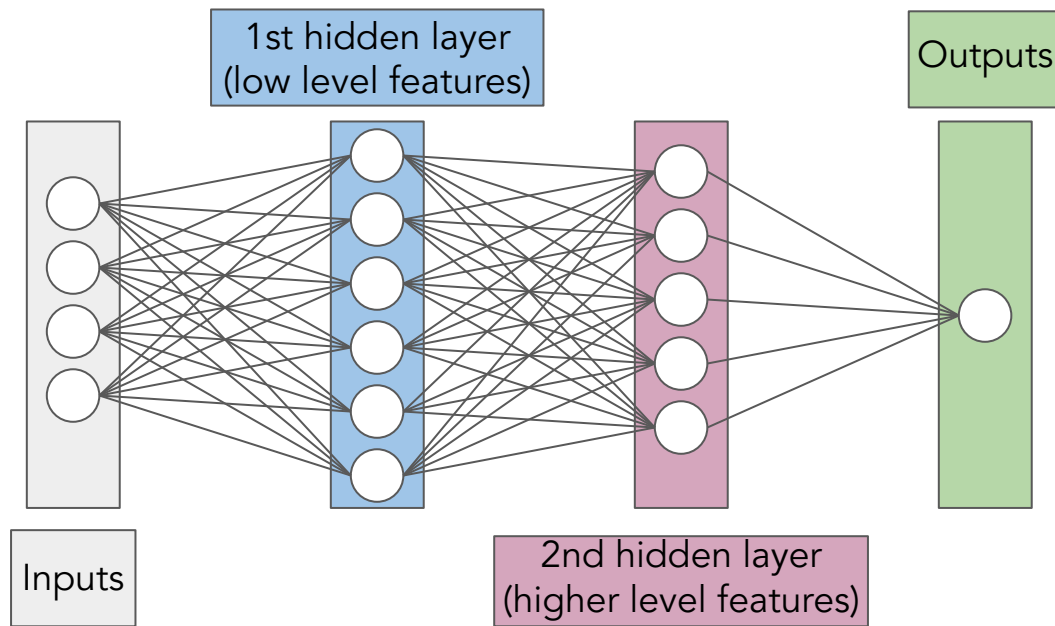


ReLU (Rectified Linear Unit):
 $\text{ReLU}(x) = \max(0, x)$



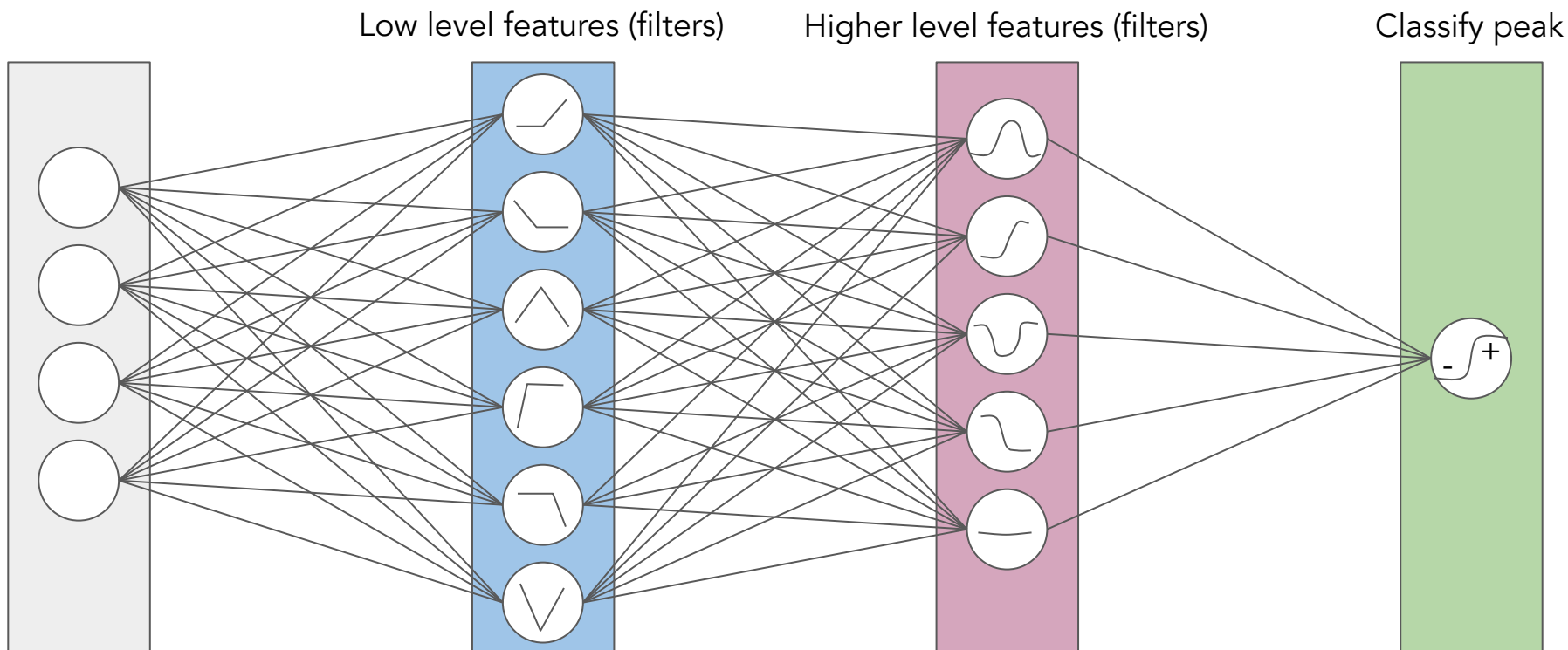
A “deep” neural network

- Each node in a hidden layer can be thought of as learning a specific feature
- Nodes closer to the inputs learn lower level features
- Nodes closer to the outputs learn more complex features



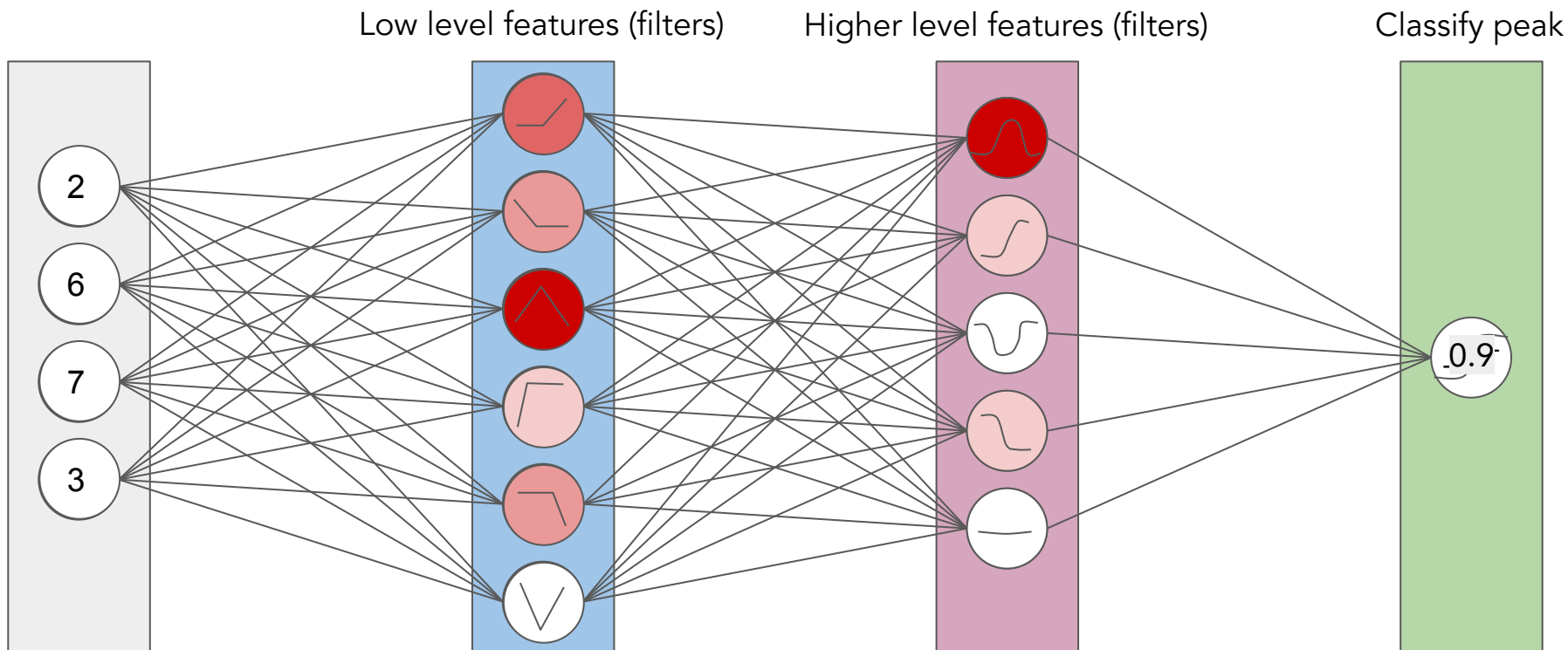
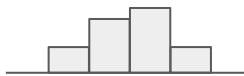
Features (filters)

Let's say our inputs are signal intensity across a window of base positions and we want to classify peaks



Features (filters)

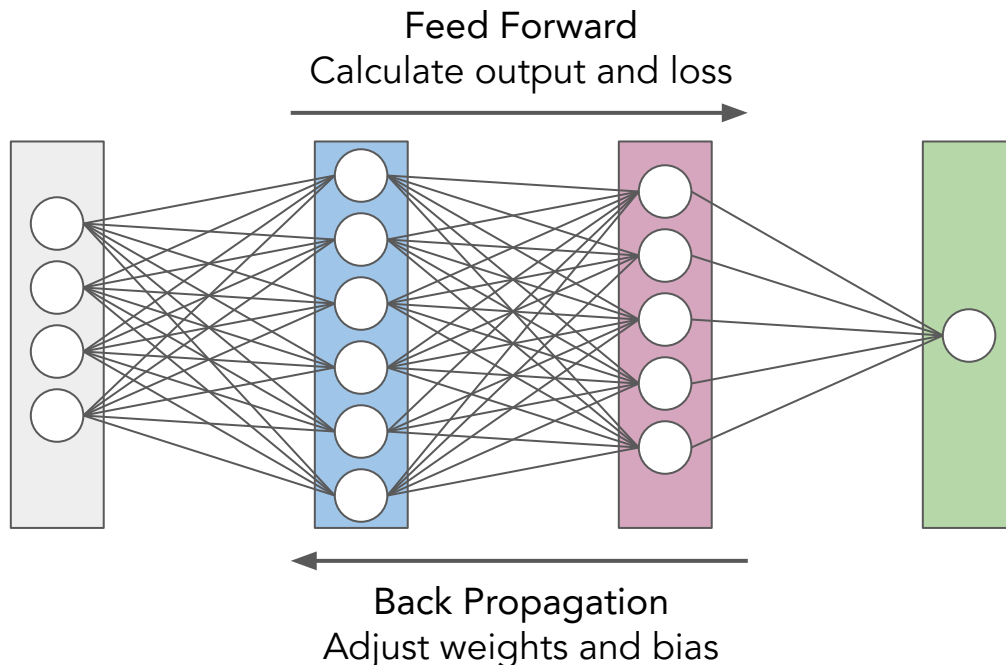
Given the input signal



Training a neural network

- Requires a large high-quality training set of ground-truth labels/values
- Need to define a loss function (how well is each prediction?)

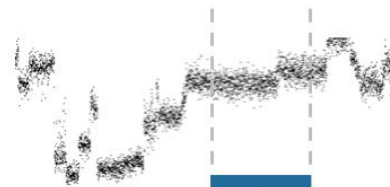
Information only flows forward
(each node in a layer is independent of every other node in that layer)



Guppy - a recurrent neural network basecaller

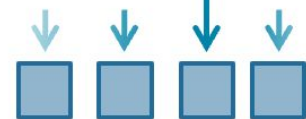
- Different low-level features are captured in the first layer across multiple windows
- Higher level features are compiled from low-level features
- Information is shared between high-level features
- Base predictions (probabilities) output for multiple positions
- All probabilities for a given base position contribute to determining base call

Base calling (RNN, raw)

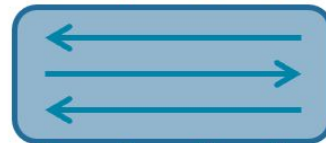


Parameters learned from training data

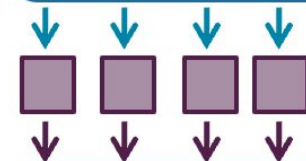
Extraction of blocks of features



Bidirectional information flow



Multi-base prediction

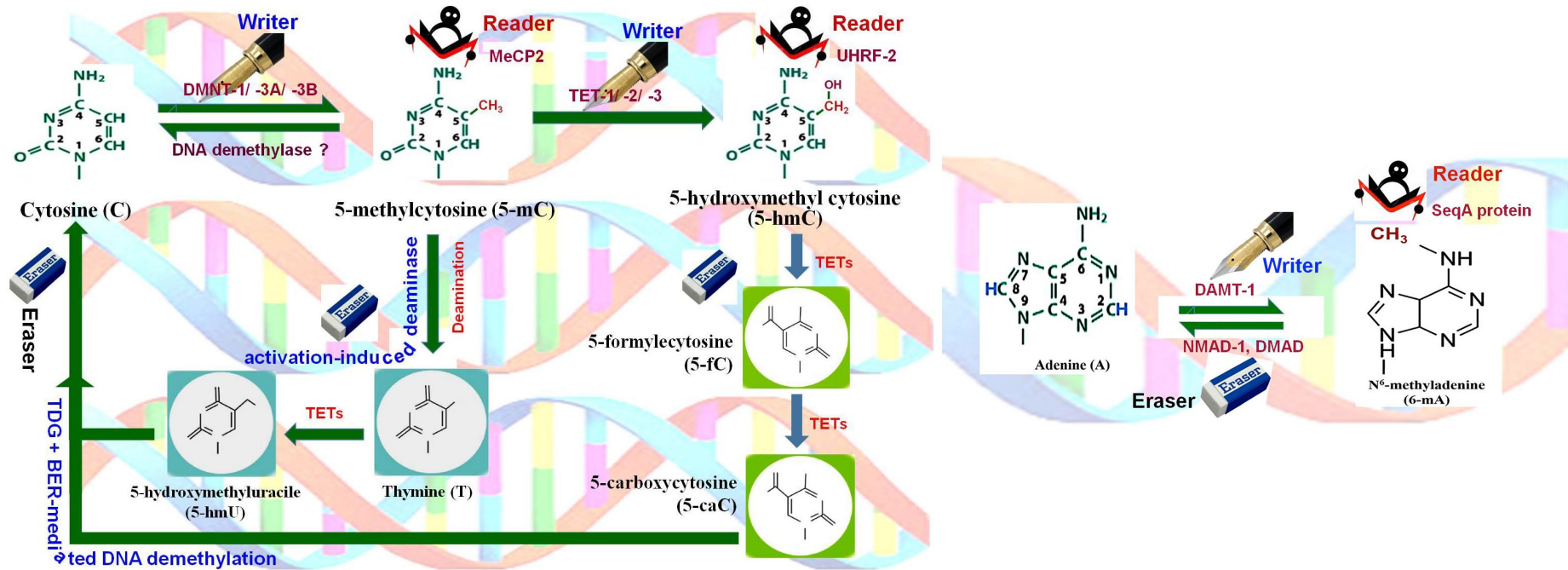


Decode to sequence



DNA/RNA base modification

Although DNA and RNA have a simple 4 letter alphabet, both have a number of possible chemical modifications that impact regulation and function



5mC



- The most common DNA base modification is methylation of the fifth carbon of cytosine (5mC)
- ~1% of the human genome is composed of 5mC
- 5mC has multiple functions
 - Silencing retroviral elements
 - Tissue-specific gene regulation
 - X-inactivation
 - Genomic imprinting

Suppression of Transposable Elements (TEs)

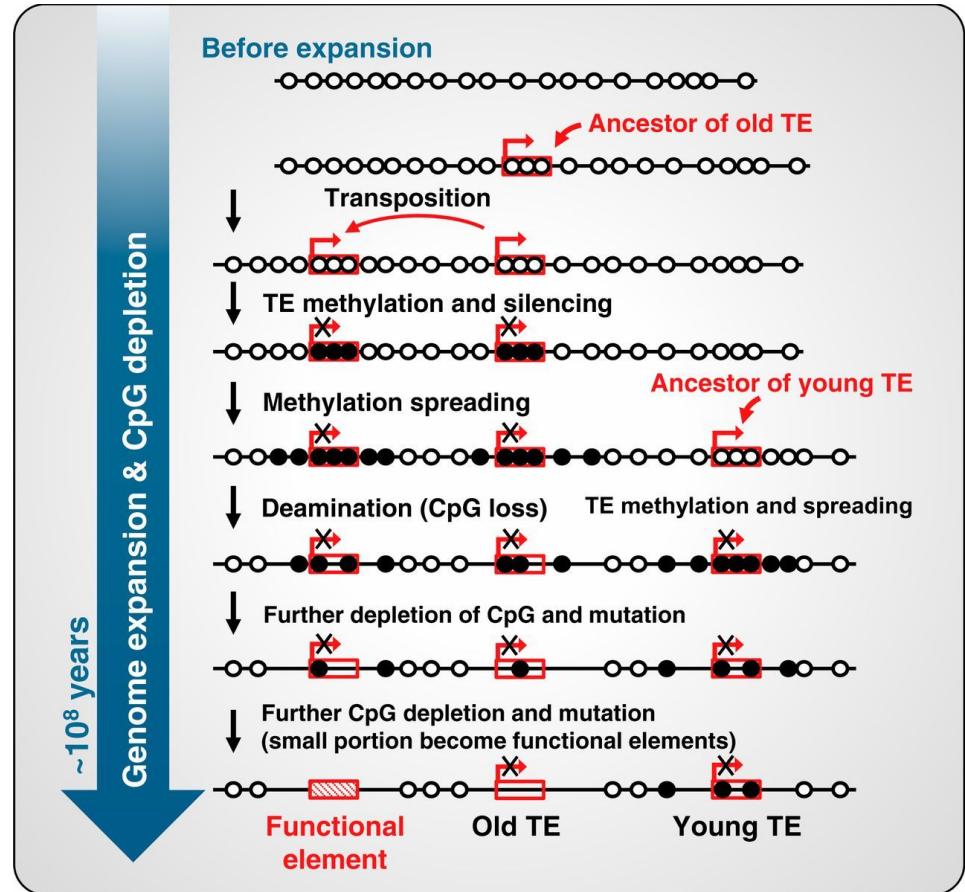
TEs are mobile repetitive sequences that can replicate and integrate into the genome

Two classes:

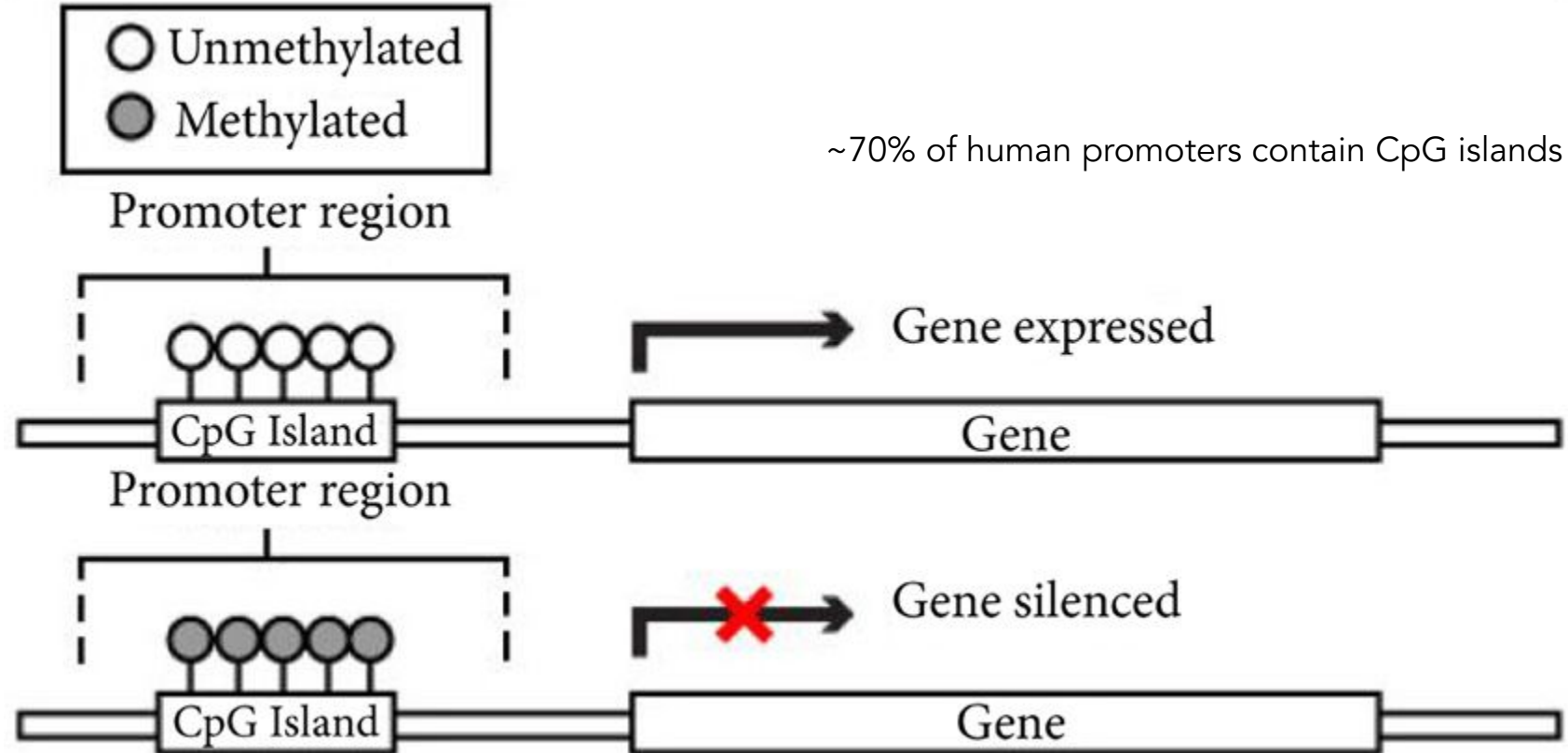
Class 1 - RNA intermediate

Class 2 - DNA intermediate

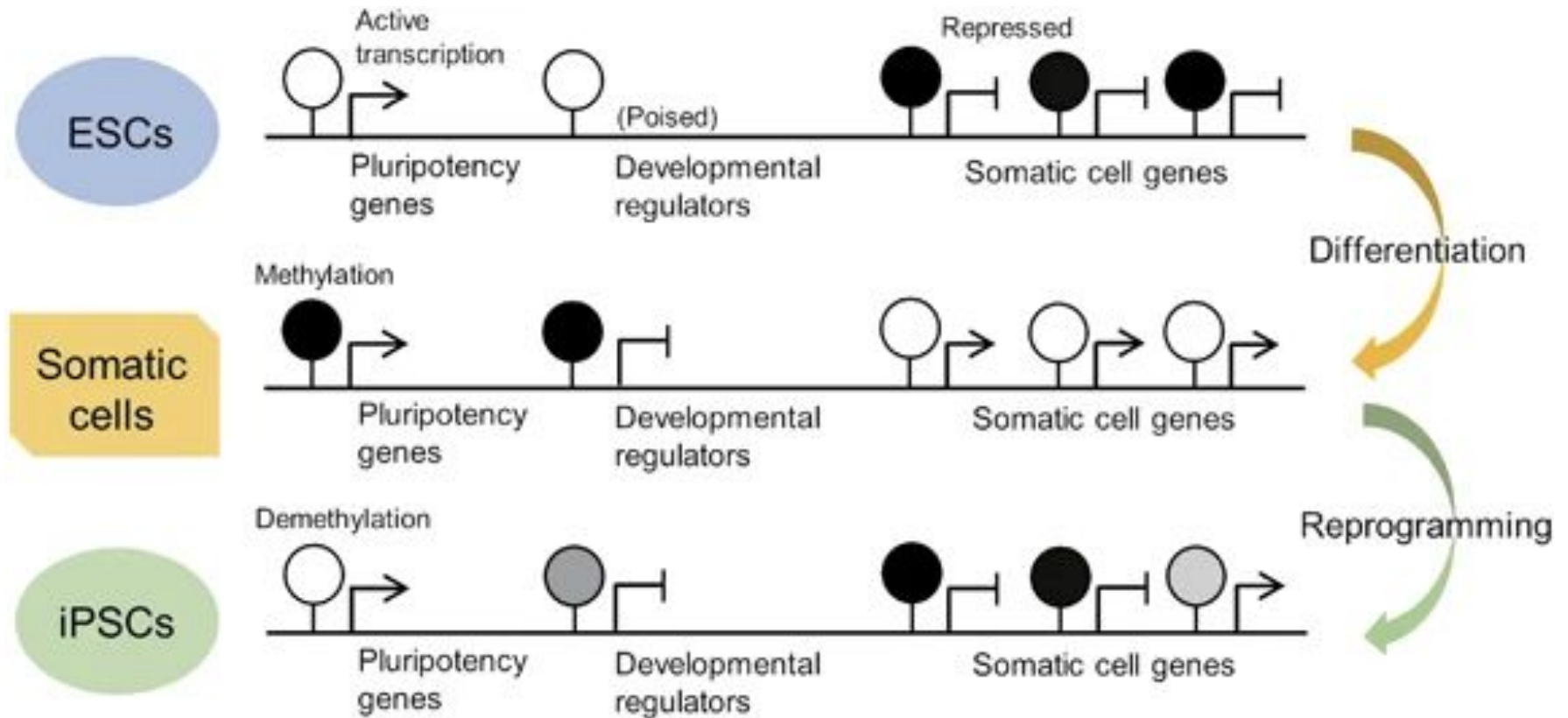
Make up ~45% of the human genome



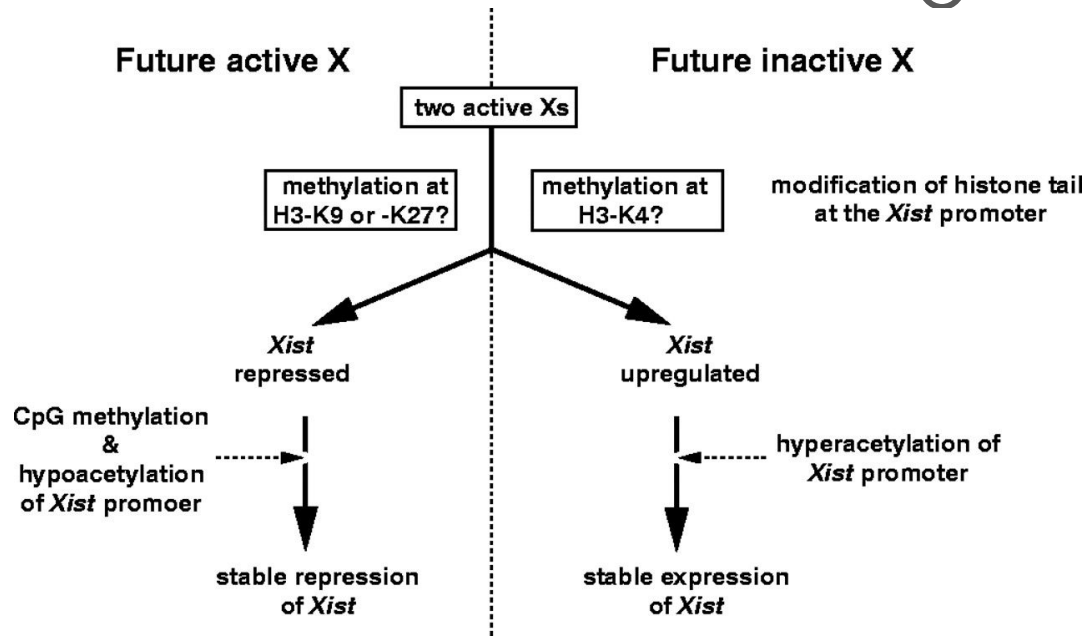
CpG Methylation



Methylation reinforces differentiation

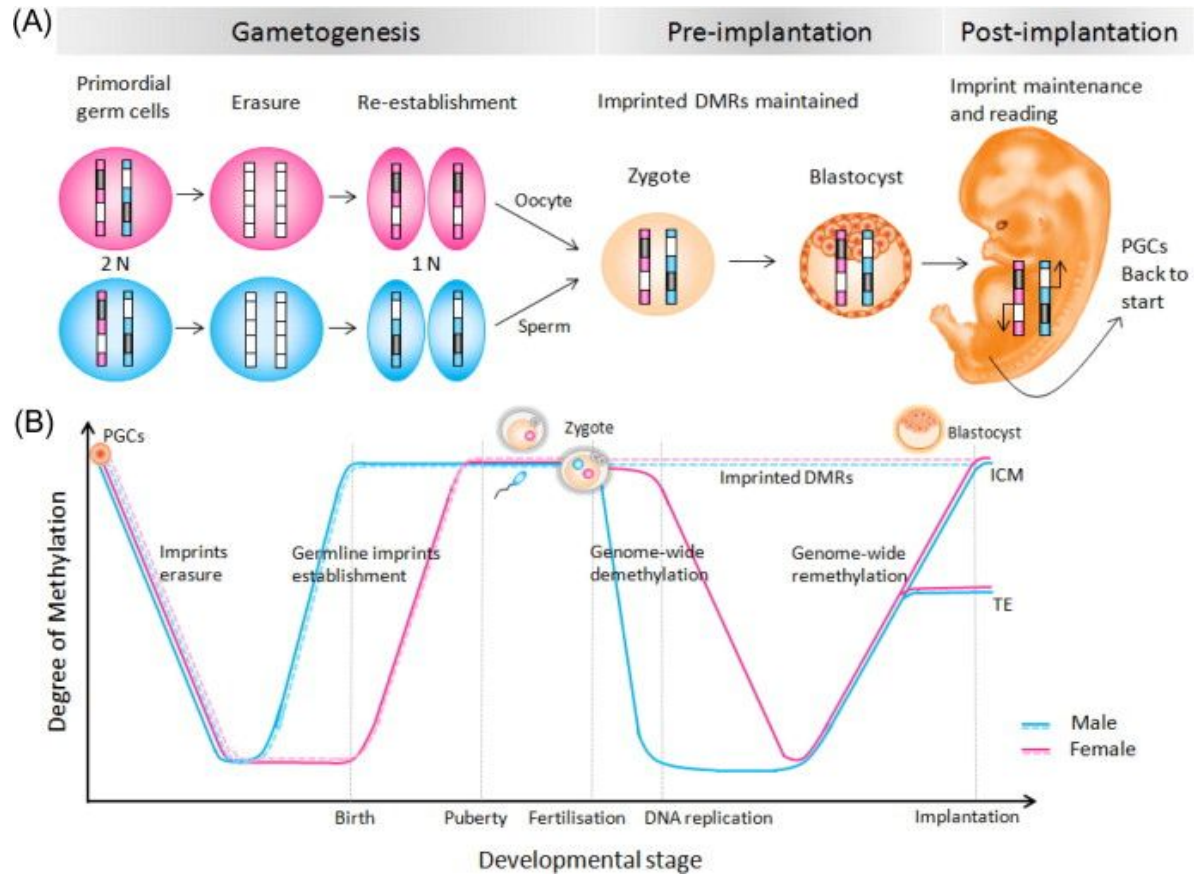


X Chromosome Inactivation



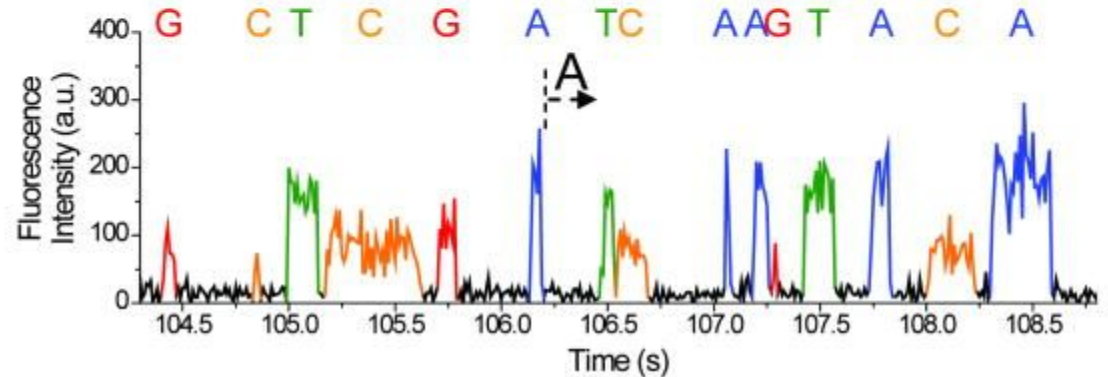
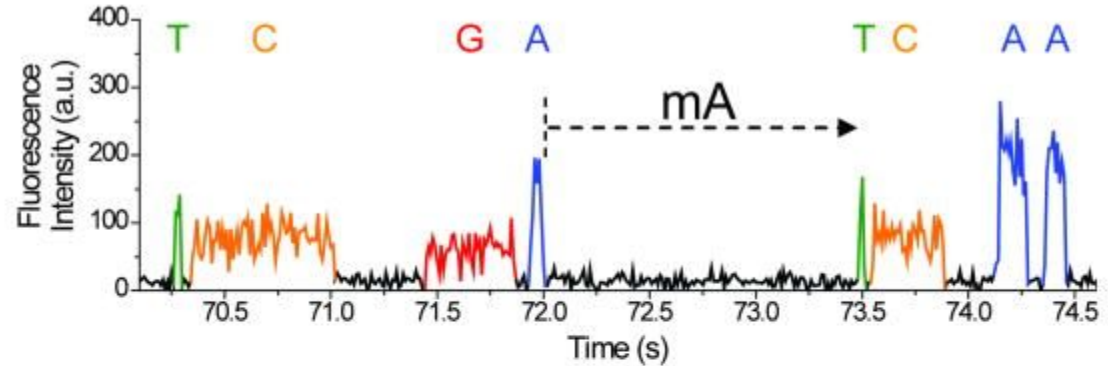
Sado T, Okano M, Li E, Sasaki H. De novo DNA methylation is dispensable for the initiation and propagation of X chromosome inactivation. *Development*. 2004 Mar;131(5):975-82.

Imprinting



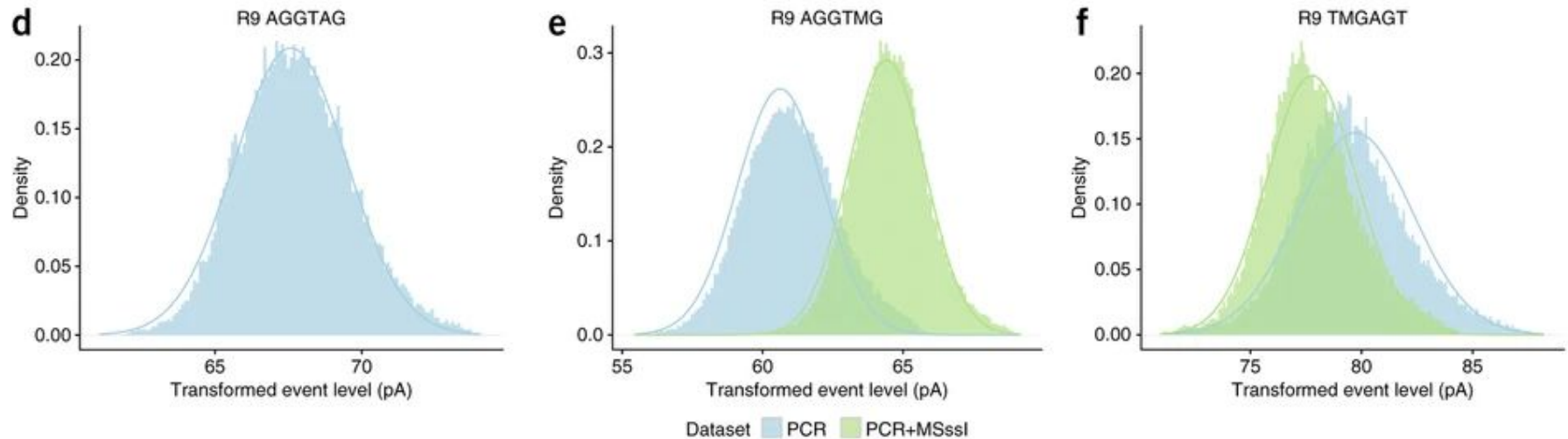
Direct sequencing of DNA methylation

PacBio



Direct sequencing of DNA methylation

Nanopore



M = 5mC

MSssl = CpG Methyltransferase