

# Introdução ao uso de dados genômicos em Sistemática

Sandra Reinales & Luana J. Sauthier  
[spreinalesl@gmail.com](mailto:spreinalesl@gmail.com); [sauthier@ib.usp.br](mailto:sauthier@ib.usp.br)

Departamento de Botânica, Universidade de São Paulo



2024



# Plano de aula

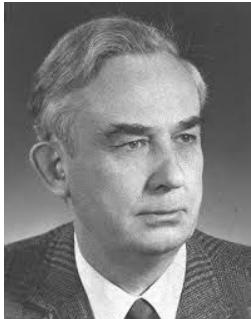
- Contextualização
- Técnicas de sequenciamento: Sanger vs. HTS
- Panorama geral dos métodos mais usados na sistemática hoje:
  - RAD-Seq
  - Genome skimming
  - Target capture
- Fluxo de trabalho no target capture
- Considerações finais
- Para saber mais



# Contextualização

## Fontes de dados para inferência filogenética: MORFOLOGIA

### Inícios da sistemática



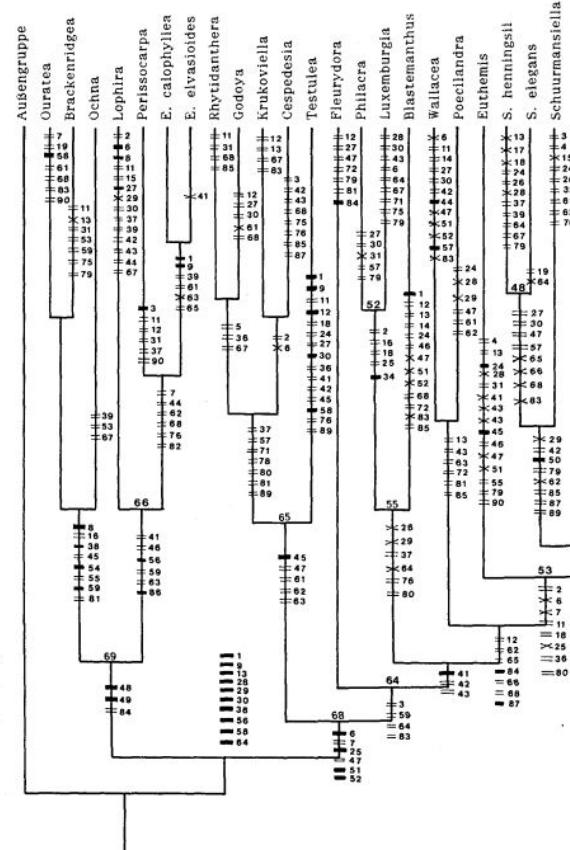
Cladística  
Willi Hennig (1950-1966)



Modelos probabilísticos  
Joseph Felsenstein (1981)



Modelos probabilísticos  
Anthony Edwards & Lucas Cavalli-Sforza (1964)



= Apomorphe Merkmalszustände, die nur einmal im Kladogramm auftreten.

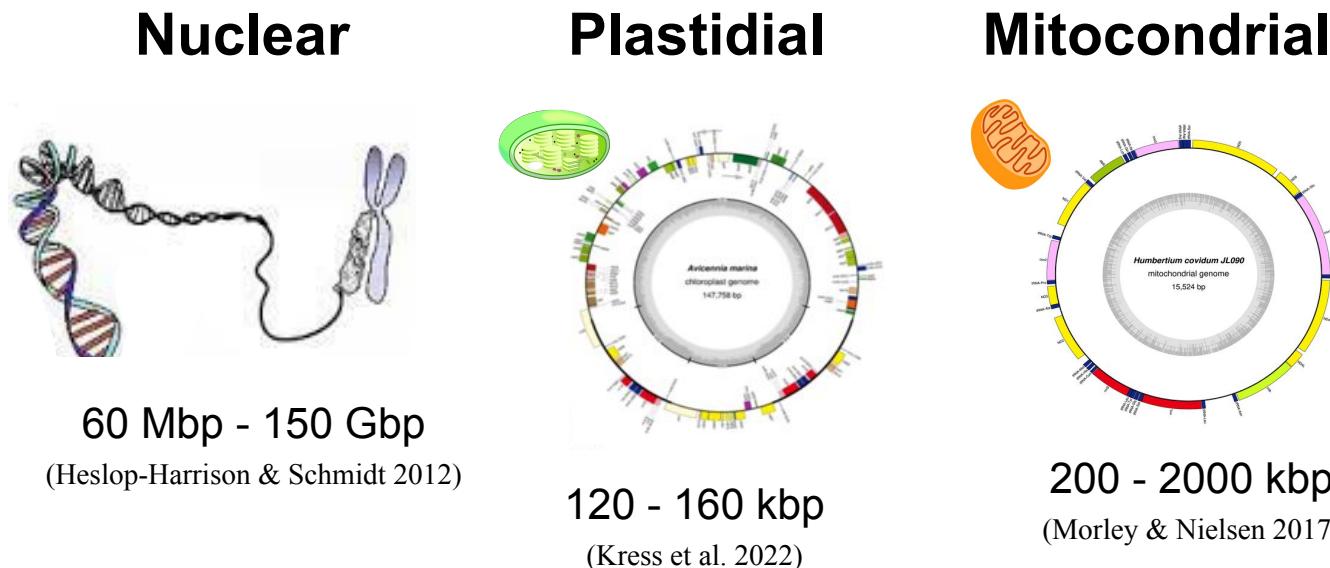
= Innerhalb des Kladogramms mehrfach auftretende Merkmalszustände (Parallelismen).

✗ = Umkehrungen von Merkmalszuständen ("reversals").

# Contextualização

## Fontes de dados para inferência filogenética: DNA

- Sanger sequencing (d. 1977 - c. 1986)
  - *rbcL* (inícios 90's)
  - ITS, *ndhF*, *matK* (1995)

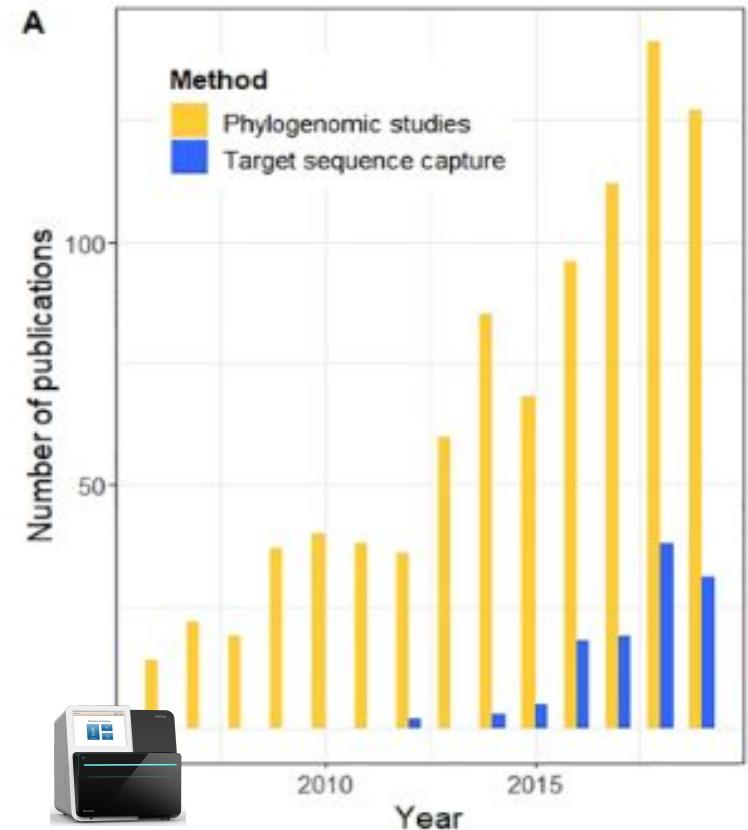


O tamanho dos genomas varia consideravelmente nas plantas

# Contextualização

Fontes de dados para inferência filogenética: DNA

- Desenvolvimento do NGS iniciou na década dos 90's
- Solexa - Cambridge Chemistry Department (1998)
- Primeiro genoma seq. (2005)
- Genome Analyzer (2006)
- Solexa - Illumina (2007)



Andermann *et al.* Front. Genet. 10 (2020)

# Contextualização

Fontes de dados para inferência filogenética: DNA

Next Generation Sequencing (NGS)

ou

High Throughput Sequencing (HTS)

illumina®

 Roche  
454  
SEQUENCING

ion torrent



PacBio

Oxford  
**NANOPORE**  
Technologies

# Contextualização

Fontes de dados para inferência filogenética: DNA

## 2º Geração

Short reads

DNA fragmentado

illumina®



ion torrent



- Genome skimming
- RAD-Seq
- Target sequence capture

## 3º Geração

Long reads

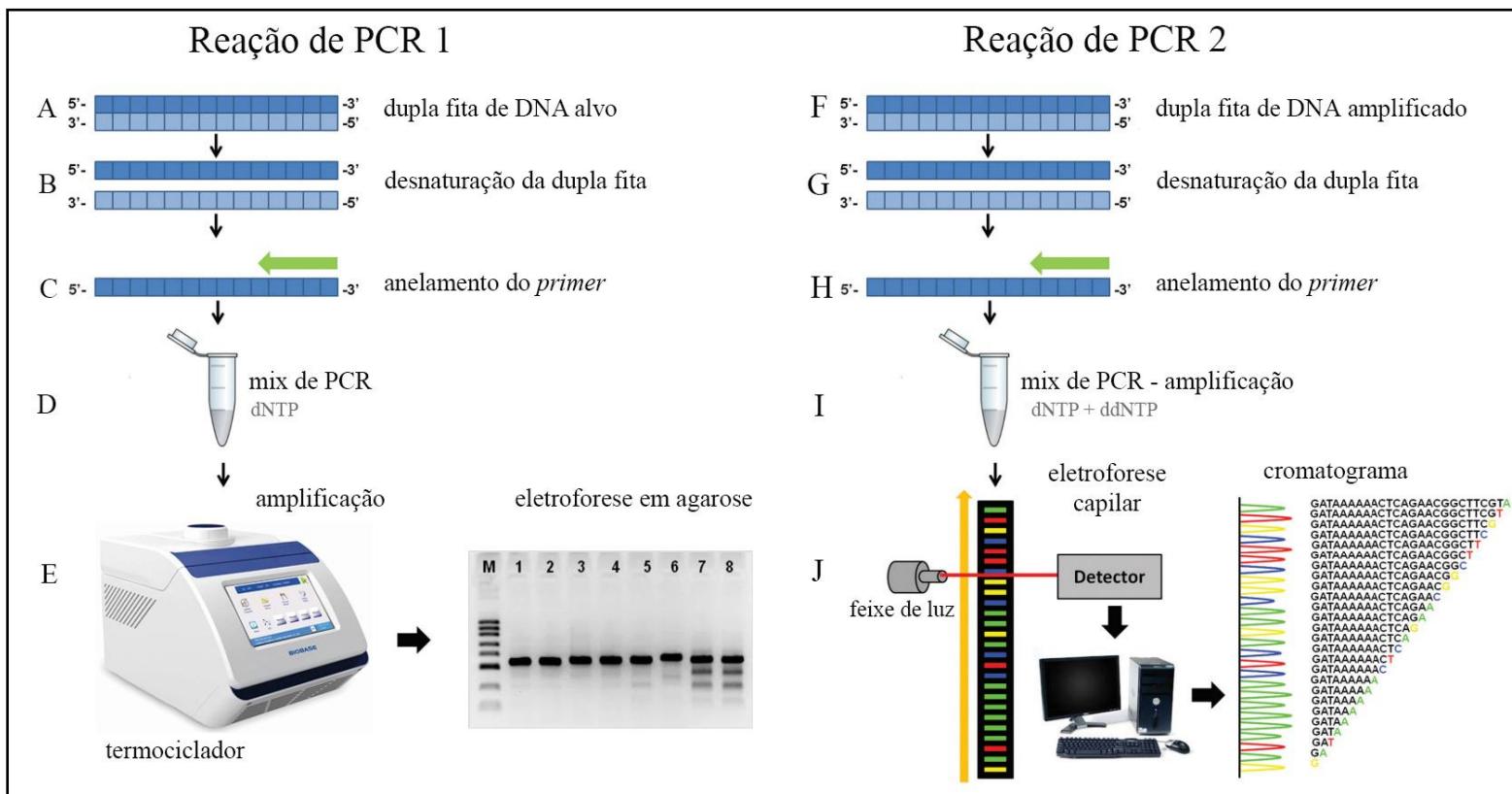
DNA não fragmentado

Oxford  
**NANOPORE**  
Technologies

PacBio

# Técnicas de sequenciamento

Sanger (1977)

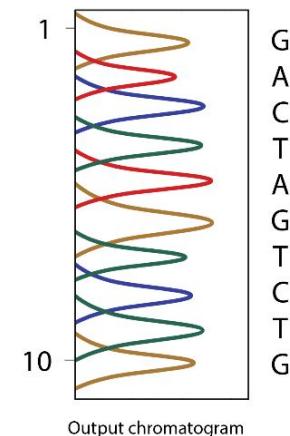
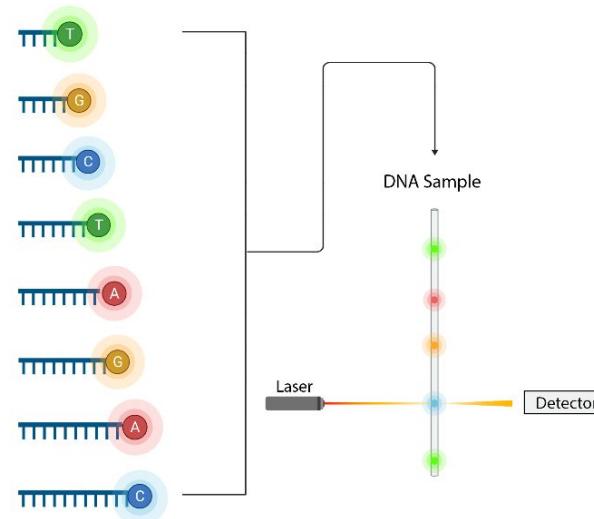
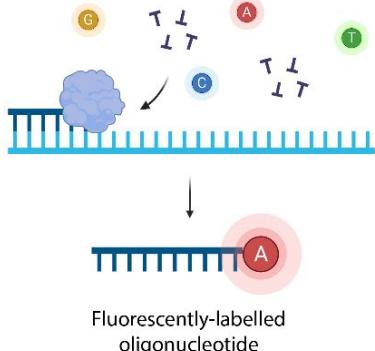
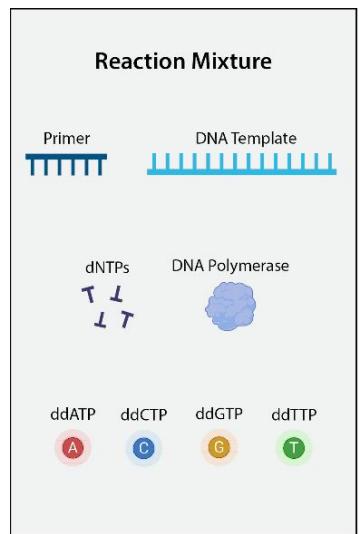


**Figura 1:** Sequenciamento Sanger diferenciando as etapas um (A-E) e dois (F-J) de reações de PCR . (A-B; F-J) desnaturação da dupla fita do DNA alvo; (C; G) anelamento dos *primers*; (D-E; I-J) transcrição da nova fita; (J) sequenciamento. Fonte: adaptado de Fietto & Maciel (2015).

# Técnicas de sequenciamento

Sanger (1977)

## Sanger Sequencing



1 Chain-termination PCR using fluorescent ddNTPs

2 Size separation and sequence analysis using capillary gel electrophoresis and fluorescence detection

# Técnicas de sequenciamento

## Comparação entre tecnologias de sequenciamento

	Sanger	NGS short reads
Input (amostra)	Clones de uma região específica do DNA (produtos de PCR)	Biblioteca de DNA do genoma completo ou regiões selecionadas
Número de amostras	96-384	1-384 por canaleta (8 canaletas)
Duração da preparação das amostras para o sequenciamento	1 dia	3-4 dias
Quantidade de dados gerados	1 fragmento (2 seqs) /amostra	Milhões de sequencias (reads) /amostra
Tamanho do fragmento/sequencia	300-1500 bp	100-500 bp
Processamento dos dados	Simples, rápido, não requer programação	Mais complexa, algumas horas/semanas requer HPC e programação

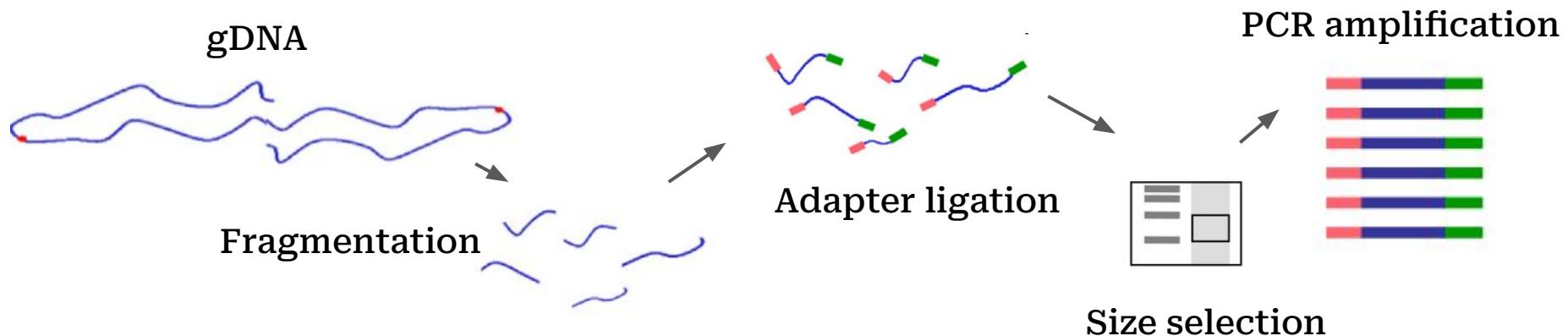
Adaptado de: <https://github.com/carol-siniscalchi>

# O que é uma biblioteca genômica?

## Definição

Uma biblioteca para NGS é uma coleção de fragmentos de DNA de tamanhos semelhantes que contém sequências adaptadoras conhecidas nas extremidades 5' e 3'.

- Uma biblioteca corresponde a uma única amostra.
- Podem ser agrupadas e sequenciadas na mesma corrida.



# O que é uma biblioteca genômica?

Kits comerciais

NEBNext® Ultra™ II for DNA Library Prep



NexTera XT DNA Library Preparation Kit



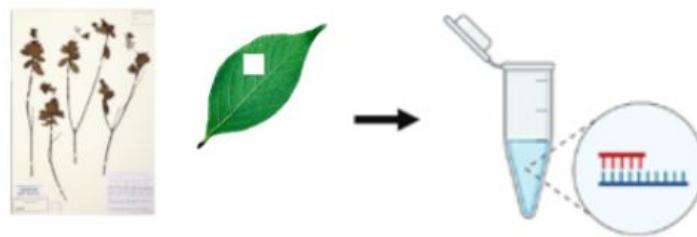
KAPA HyperPrep Kits



# O que é uma biblioteca genômica?

Etapas: extração do DNA

## (A) Extração de DNA com o protocolo CTAB

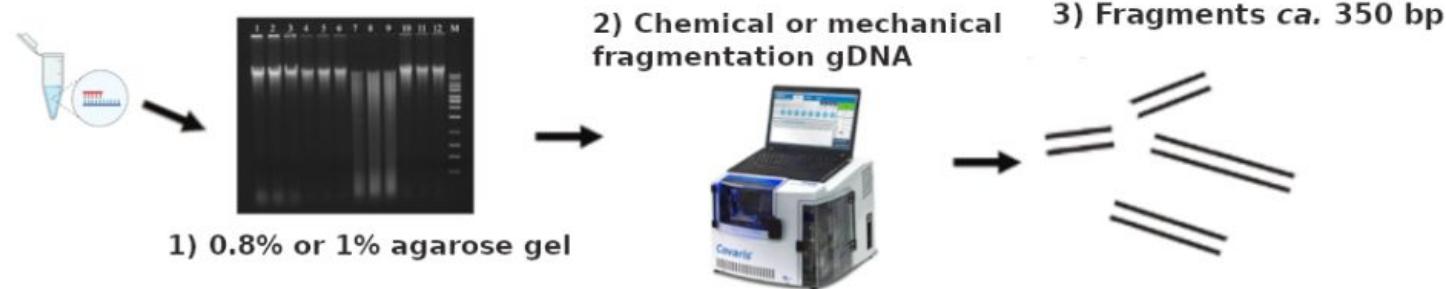


RECOMMENDED INPUT AMOUNTS	TOTAL DNA
Ultra II DNA Library Prep Kit (NEB #E7645)*	500 pg–1 µg

# O que é uma biblioteca genômica?

Etapas: análise do DNA extraído

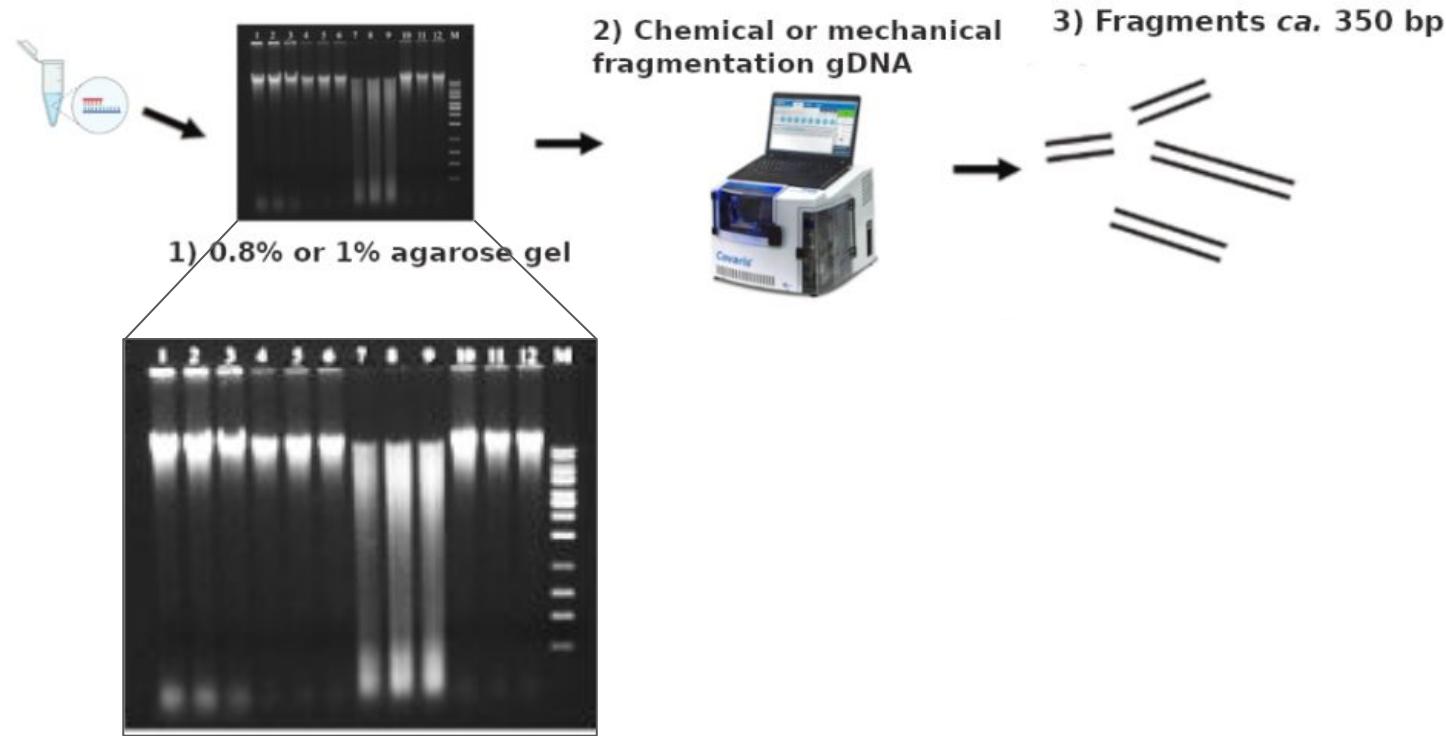
## (B) Análise do tamanho dos fragmentos e fragmentação do DNA



# O que é uma biblioteca genômica?

Etapas: análise do DNA extraído

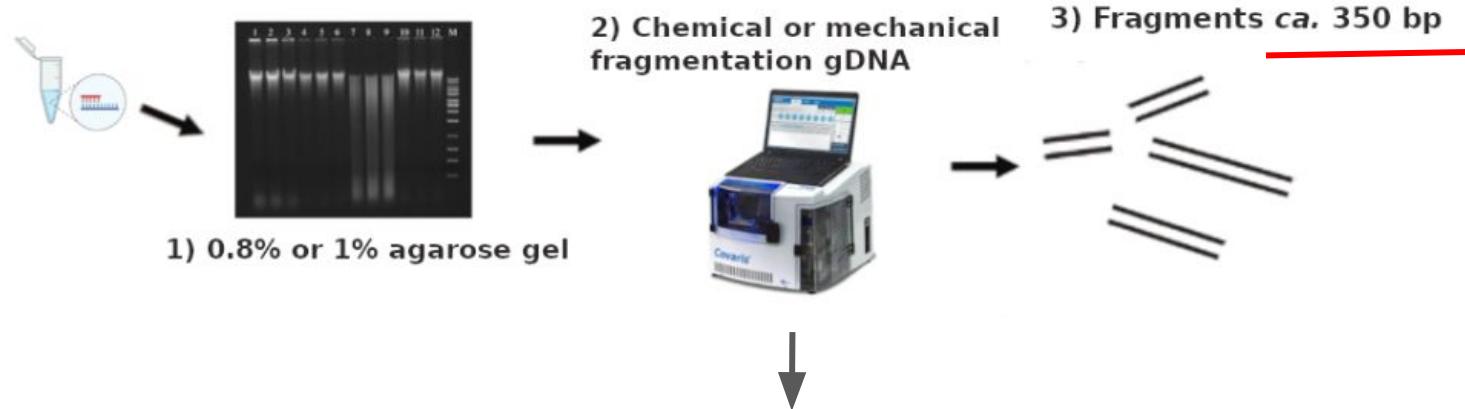
## (B) Análise do tamanho dos fragmentos e fragmentação do DNA



# O que é uma biblioteca genômica?

Etapas: fragmentação do DNA

## (B) Análise do tamanho dos fragmentos e fragmentação do DNA

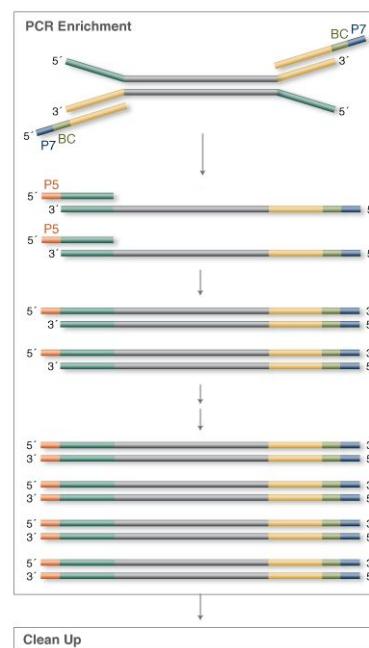
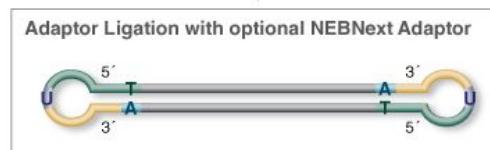
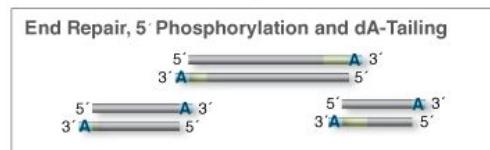
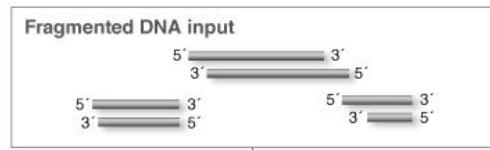
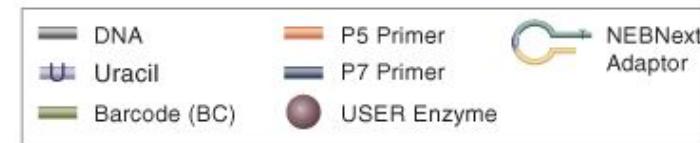


- Covaris - sonicador UDS \$5.28 c/u
- Enzimas de restrição

# O que é uma biblioteca genômica?

## Etapas: ligação de adaptadores

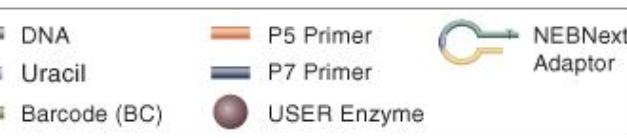
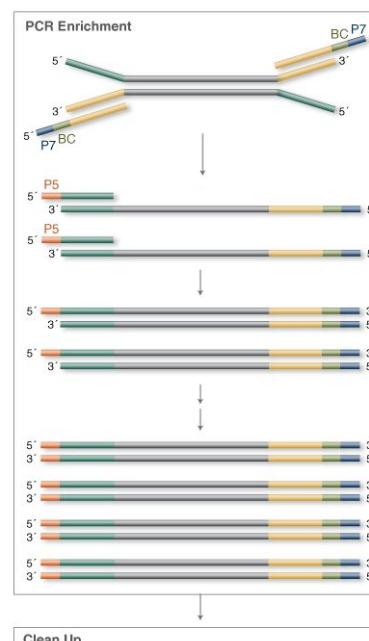
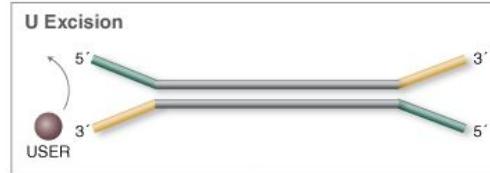
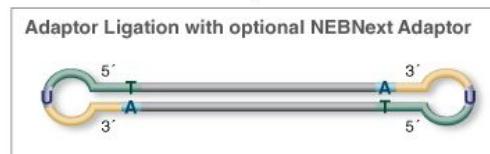
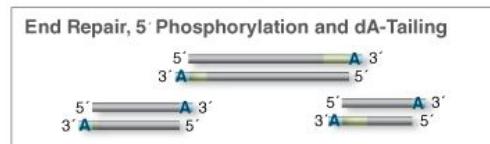
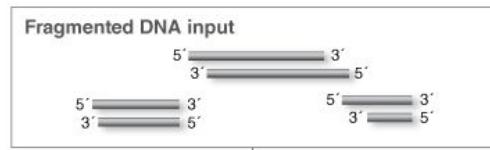
### (C) Preparação das bibliotecas com kit



# O que é uma biblioteca genômica?

Etapas: ligação de adaptadores

## (C) Preparação das bibliotecas com kit



P5 Primer  
**AATGATAACGGCGACCACCGAGATCTACAC**

P7 Primer  
**ATCTCGTATGCCGTCTTCTGCTTG**

# O que é uma biblioteca genômica?

Etapas: ligação de adaptadores

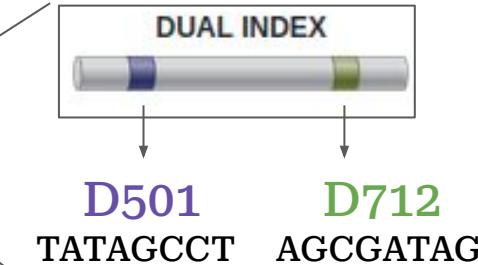
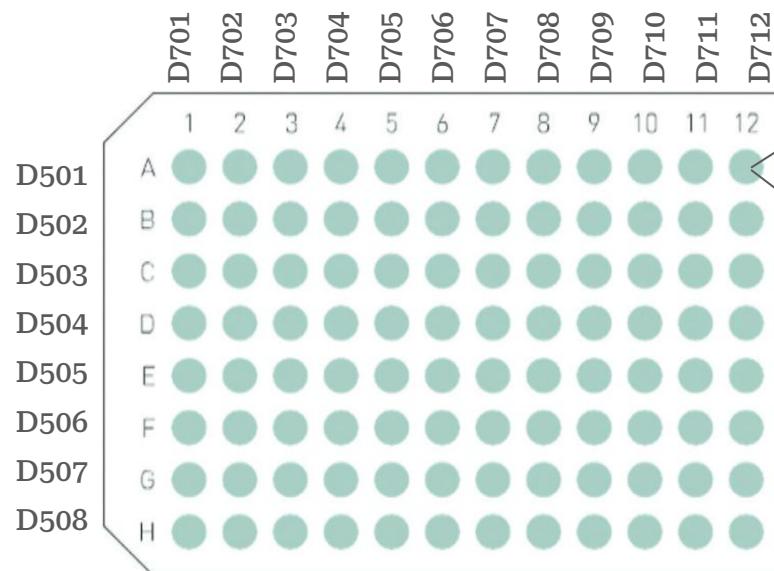
## (C) Preparação das bibliotecas com kit



NEBNext® Multiplex Oligos Selection Chart

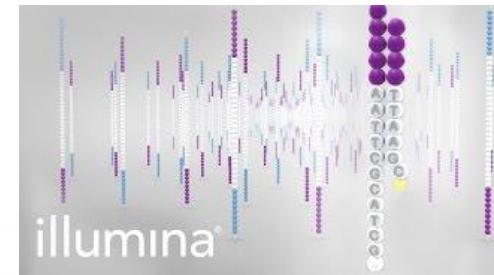


Multiplexing



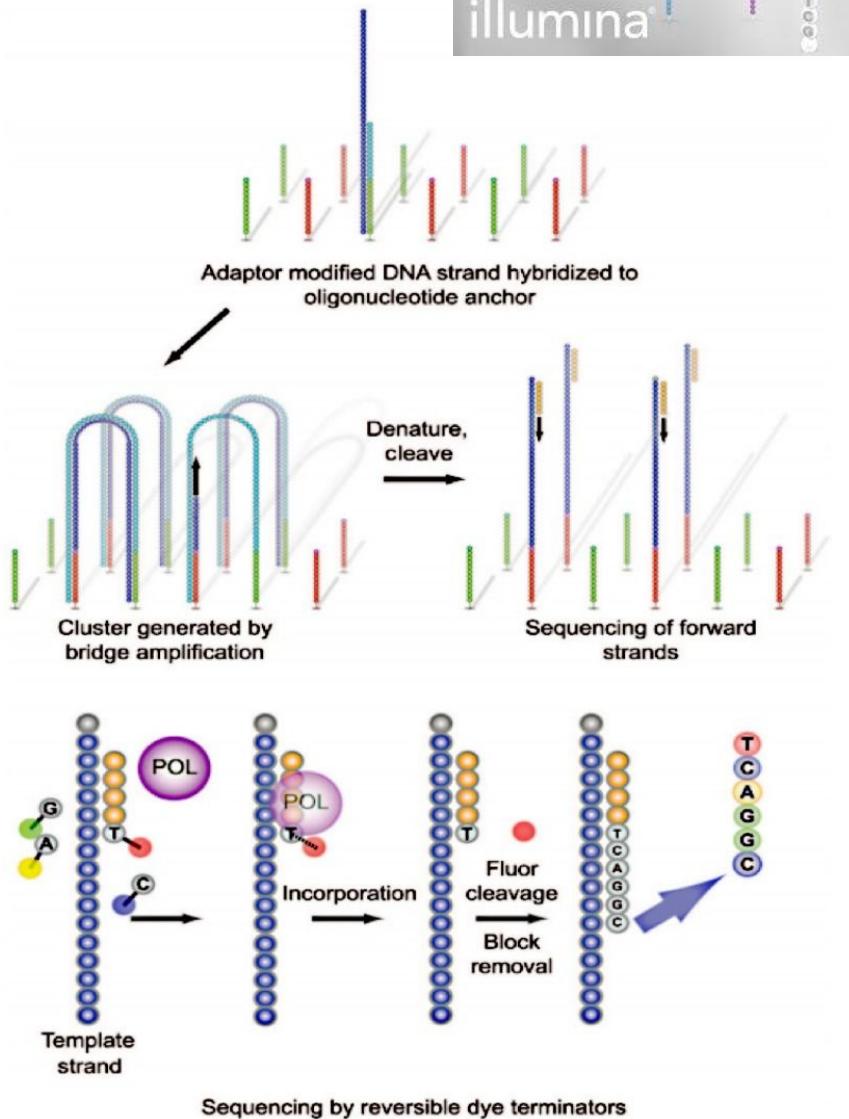
# Técnicas de sequenciamento

## Sequenciamento Illumina



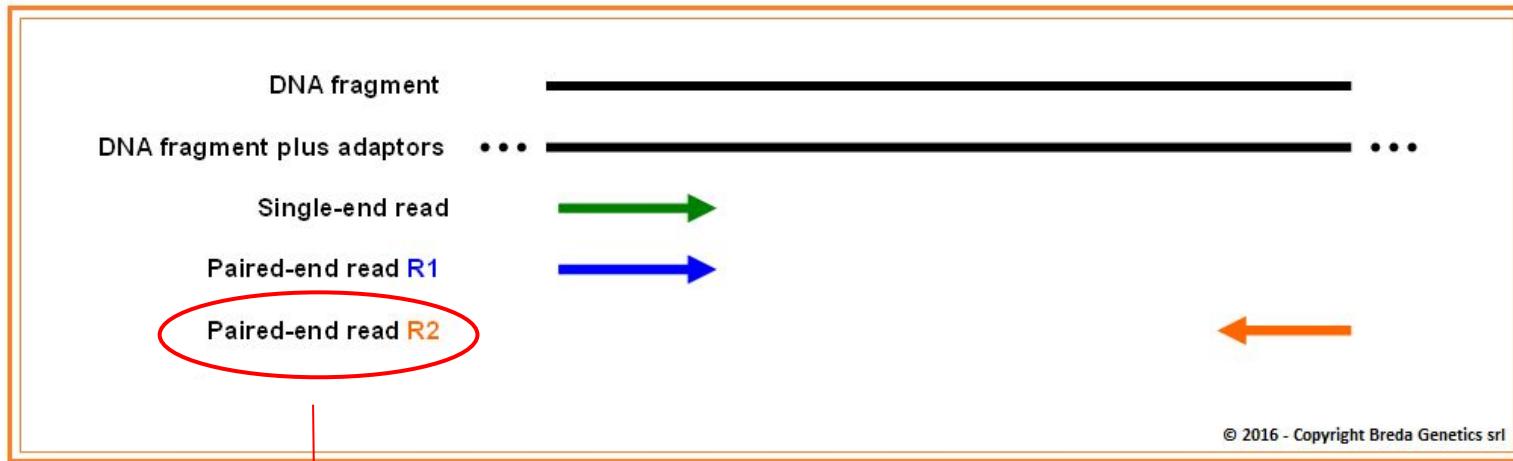
Três etapas principais:

1. Ligação das bibliotecas na plataforma
2. PCR fase sólida (Bridge amplification)
3. Sequenciamento por síntese



# Técnicas de sequenciamento

## Sequenciamento Illumina



<https://bredagenetics.com/single-end-reads-e-paired-end-reads/?lang=it>

SRR10951301\_1.fastq.gz

SRR10951301\_2.fastq.gz

# Técnicas de sequenciamento

## Plataformas Illumina

MiniSeq



MiSeq



NextSeq



HiSeq 4000



HiSeq X Ten



MAX OUTPUT  
**8 Gb**  
MAX READ NUMBER  
**25 million**  
MAX READ LENGTH  
**2x150 bp**

MAX OUTPUT  
**15 Gb**  
MAX READ NUMBER  
**25 million**  
MAX READ LENGTH  
**2x300 bp**

MAX OUTPUT  
**120 Gb**  
MAX READ NUMBER  
**400 million**  
MAX READ LENGTH  
**2x150 bp**

MAX OUTPUT  
**1500 Gb**  
MAX READ NUMBER  
**5 billion**  
MAX READ LENGTH  
**2x150 bp**

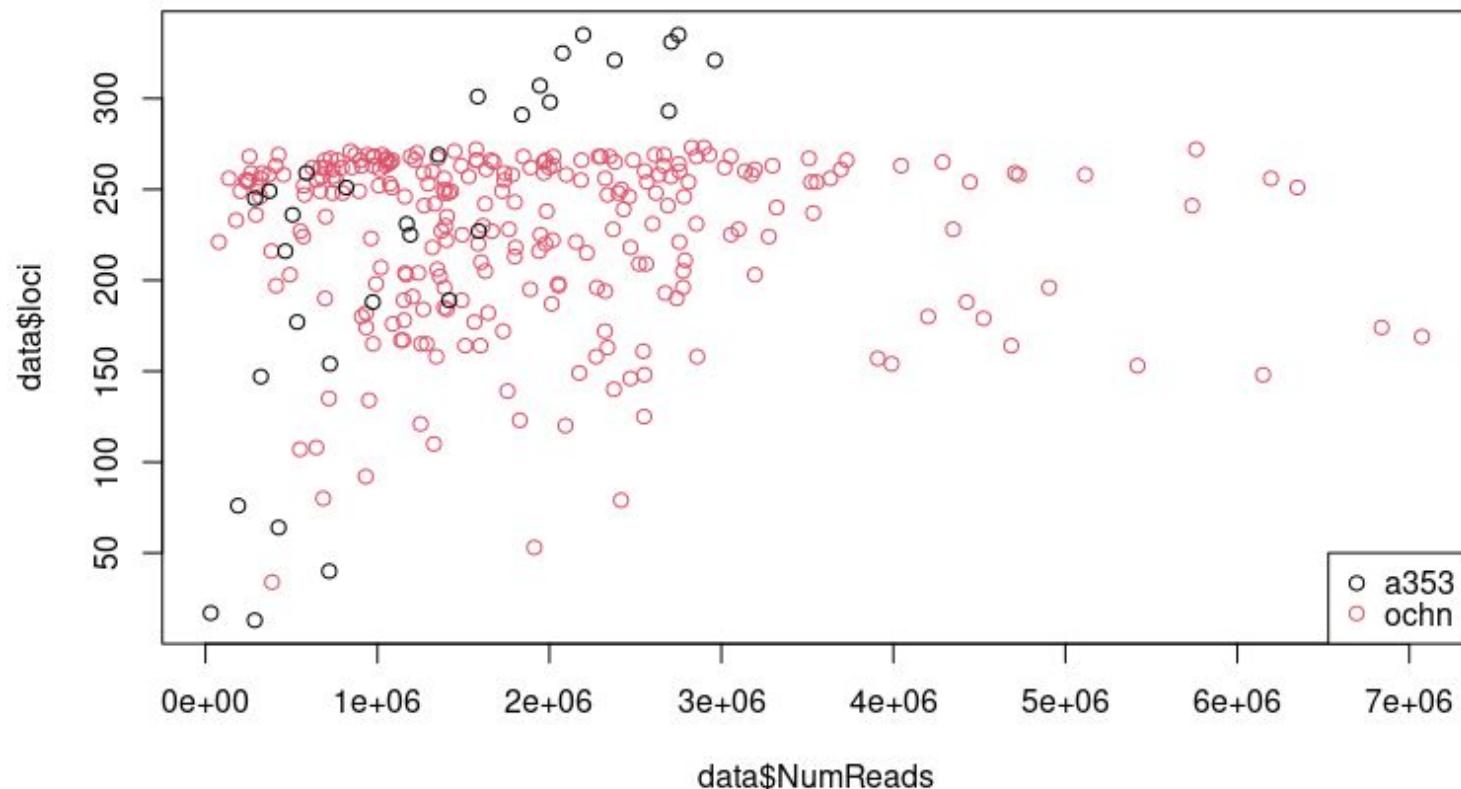
MAX OUTPUT  
**1800 Gb**  
MAX READ NUMBER  
**6 billion**  
MAX READ LENGTH  
**2x150 bp**

<https://www.illumina.com/systems/sequencing-platforms.html>

6 bilhões / 8 canaletas = 750 milhões reads → 750 milhões / 96 amostras = ~ 8 milhões de reads por amostra

# Técnicas de sequenciamento

## Sequenciamento Illumina



# **Os dados chegaram, e agora?**



# O que é um pipeline?

Conjunto de ferramentas (pacotes, programas e scripts):

- Instalação conjunta
- Estão interligadas: saída de uma é a entrada da outra
- Facilita análises padrão



# O que é um pipeline?

Conjunto de ferramentas (pacotes, programas e scripts):

- Instalação conjunta
  - Estão interligadas: saída de uma é a entrada da outra
  - Facilita análises padrão
- 
- Pouco flexíveis: análises mais complexas precisam de programação
  - Desconhecimento do procedimento



# O que é um pipeline?

	Read cleaning	Sequence engineering	Intron recovery	MSA generation	Allele phasing	SNP extraction	Ease of installation
Hybpiper (Johnson et al., 2016)	○	●	●	○	○	○	■
Phyluce (Faircloth, 2016)	■	○	○	●	●	■	●
Secapr (Andermann et al., 2018)	●	■	■	●	●	■	●
Hybphaser (Nauheimer, et al., 2021)	○	●	●	●	●	●	●
Captus (Ortiz, et al., 2023)	●	●	●	●	○	○	●

Adaptado de : Andermann *et al.* Front. Genet. 10 (2020)

# Comparação dos métodos

## Genome skimming, RAD-Seq, Target capture

	Genome skimming	RAD-seq	Hyb-seq
Demand for plant materials	Fresh, silica-gel dried plant tissues and specimen	Fresh, silica-gel dried plant tissues	Fresh, silica-gel dried plant tissues and specimen
Demand for DNA template quality	Low	Medium	Low
Applicable to specimen	Yes	No	Yes
Material for sequencing	Total genomic DNA	Restriction fragments	Captured loci using probes
Genome data obtained	Complete plastid genome, nrDNA, partial mitochondrial genome, coding and non-coding genes	Loci with single nucleotide polymorphism (SNP) mainly from nuclear genome; coding and non-coding genes	Targeted nuclear, plastid and/or mitochondrial loci; coding and non-coding genes
Targeted loci sequenced	Yes	No	Yes
Identification of orthologs	Easy	Difficult	Easy
Missing data among species	No	Yes	No
Taxonomic levels for phylogenetic relationships	All levels from shallow to deep	Shallow levels, below inter-generic	All levels from shallow to deep, above intraspecific

# Comparação dos métodos

## Genome skimming

Nevill et al. *Plant Methods* (2020) 16:1  
<https://doi.org/10.1186/s13007-019-0534-5>

### METHODOLOGY

Large scale genome skimming from herbarium material for accurate plant identification and phylogenomics

Paul G. Nevill<sup>1,2,3\*</sup>, Xiao Zhong<sup>4,5</sup>, Julian Tonti-Filippini<sup>4,5</sup>, Margaret Byrne<sup>2,6,7</sup>, Michael Hislop<sup>6</sup>, Kevin Thiele<sup>2,6</sup>, Stephen van Leeuwen<sup>6</sup>, Laura M. Boykin<sup>4,5</sup> and Ian Small<sup>4,5</sup>

Chloroplast genome - 96.1% of samples

Complete or near-complete nuclear ribosomal RNA gene - 93.3% of samples

frontiers | Frontiers in Plant Science

Sections Articles Research Topics Editorial board About journal

ORIGINAL RESEARCH article

Front. Plant Sci., 03 April 2022  
Sec. Plant Systematics and Evolution  
Volume 13 - 2022 | <https://doi.org/10.3389/fpls.2022.832054>

Genome Skimming Contributes to Clarifying Species Limits in *Paris* Section *Axiparis* (Melanthiaceae)

 Yunheng Ji<sup>1,2\*</sup>, Jin Yang<sup>1,3</sup>, Jacob B. Landis<sup>4,5</sup>, Shuying Wang<sup>1,3</sup>, Lei Jin<sup>1,6</sup>, Pingxuan Xie<sup>1,6</sup>, Haiyang Liu<sup>7</sup>, Jun-Bo Yang<sup>8</sup>, Ting-Shuang Yi<sup>8</sup>

“... we employed a genome skimming approach to recover the plastid genomes (plastomes) and nuclear ribosomal DNA (nrDNA)”



Plant Diversity  
Volume 46, Issue 3, May 2024, Pages 344-352



Research paper

Deep genome skimming reveals the hybrid origin of *Pseudosasa gracilis* (Poaceae: Bambusoideae)

Xiang-Zhou Hu<sup>a b 1</sup> , Cen Guo<sup>a 1</sup> , Sheng-Yuan Qin<sup>a b</sup> , De-Zhu Li<sup>a b</sup> , Zhen-Hua Guo<sup>a b</sup>

A pipeline for assembling low copy nuclear markers from plant genome skimming data for phylogenetic use

Marcelo Reginato



“By integrating ... chloroplast genome data ... and subgenome data ... together with morphological and geographic information, we provide an empirical example of bamboo hybridization

# Comparação dos métodos RAD-seq

## scientific reports

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [scientific reports](#) > [articles](#) > [article](#)

Article | [Open access](#) | Published: 23 June 2022

### Genetic diversity of loquat (*Eriobotrya japonica*) revealed using RAD-Seq SNP markers

[Yukio Nagano](#), [Hiroaki Tashiro](#), [Sayoko Nishi](#), [Naofumi Hiehata](#), [Atsushi J. Nagano](#) & [Shinji Fukuda](#)✉

[Scientific Reports](#) 12, Article number: 10200 (2022) | [Cite this article](#)

## Systematic Entomology



Original Article

### RAD-seq phylogenomics recovers a well-resolved phylogeny of a rapid radiation of mutualistic and antagonistic yucca moths

CLIVE T. DARWELL ✉ DAVID M. RIVERS, DAVID M. ALTHOFF

First published: 02 May 2016 | <https://doi.org/10.1111/syen.12185> | Citations: 17

JOURNAL ARTICLE

### Integrating restriction site-associated DNA sequencing (RAD-seq) with morphological cladistic analysis clarifies evolutionary relationships among major species groups of bee orchids Ⓢ

Richard M Bateman ✉, Gábor Sramkó, Ovidiu Paun

*Annals of Botany*, Volume 121, Issue 1, January 2018, Pages 85–105, <https://doi.org/10.1093/aob/mcx129>

Published: 09 January 2018 Article history ▾

## scientific reports

Explore content ▾ About the journal ▾ Publish with us ▾

[nature](#) > [scientific reports](#) > [articles](#) > [article](#)

Article | [Open access](#) | Published: 23 March 2015

### Identification of domestication-related loci associated with flowering time and seed size in soybean with the RAD-seq genotyping method

[Ling Zhou](#), [Shi-Bo Wang](#), [Jianbo Jian](#), [Qing-Chun Geng](#), [Jia Wen](#), [Qijian Song](#), [Zhenzhen Wu](#), [Guang-Jun Li](#), [Yu-Qin Liu](#), [Jim M. Dunwell](#), [Jin Zhang](#), [Jian-Ying Feng](#), [Yuan Niu](#), [Li Zhang](#), [Wen-Long Ren](#) & [Yuan-Ming Zhang](#)



## Molecular Phylogenetics and Evolution

Volume 126, September 2018, Pages 1–16



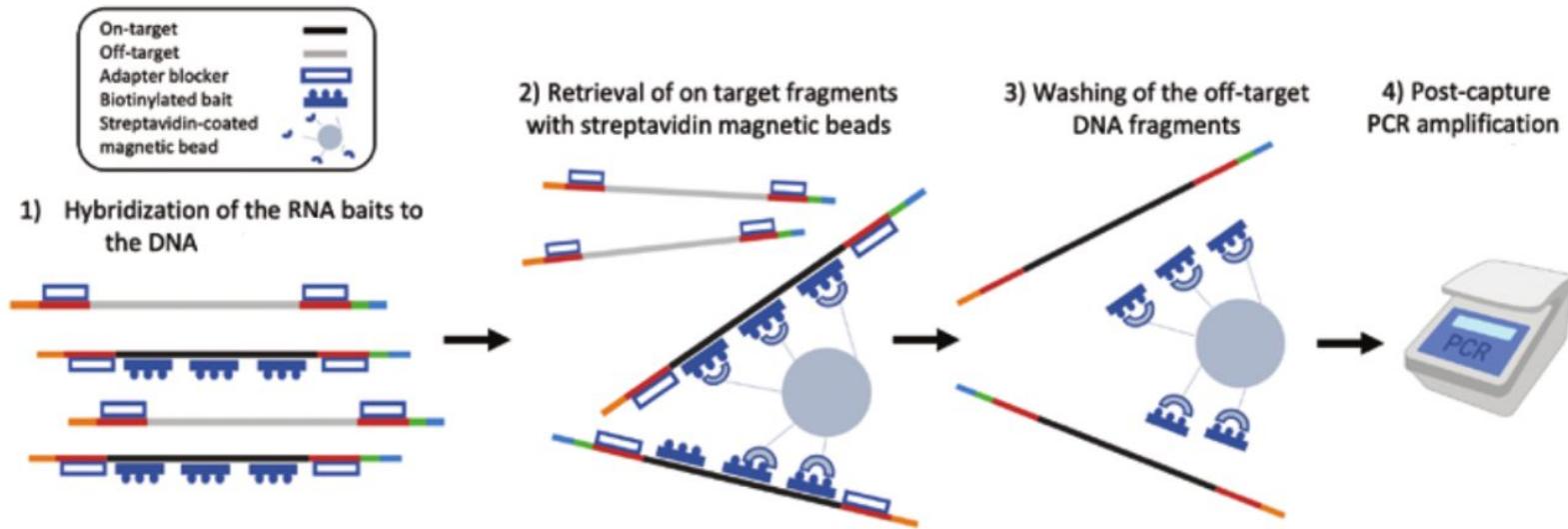
### Resolving relationships and phylogeographic history of the *Nyssa sylvatica* complex using data from RAD-seq and species distribution modeling

[Wenbin Zhou](#)<sup>a</sup>, [Xiang Ji](#)<sup>b,c</sup>, [Shihori Obata](#)<sup>a</sup>, [Andrew Pais](#)<sup>a</sup>, [Yibo Dong](#)<sup>a</sup>, [Robert Peet](#)<sup>d</sup> ✉, [Qiu-Yun \(Jenny\) Xiang](#)<sup>a</sup> ✉

# O que são as probes?

Um fragmento de oligonucleotídeos usualmente de RNA de 80-120 bp, desenhados para capturar fragmentos alvo numa biblioteca de DNA genômico.

## (E) Captura dos fragmentos de interesse com RNA *baits* mediante hibridização (enriquecimento das bibliotecas)



# O que são as probes?

- Geralmente capturam regiões de baixa cópia do genoma nuclear
- Podem ser universais ou específicas

**Angiosperms353**  
(Johnson et al. 2019)



**CompCOS**  
(Mandel et al. 2017)



**Ochnaceae**  
(Schneider et al. 2021)



# Probes universais vs. específicas

	Universais	Específicas
Poder de resolução a nível infraespecífico	Baixo*	Alto
Especificidade das probes	Baixa	Alta
Custo do desenho das probes	Baixo	Alto
Compatibilidade entre diferentes estudos	Alta	Baixa*

# Probes universais vs. específicas

Universais      Específicas

Pode ser  
a nível de:

Específico

Custodiado

Composto  
diferente



PROTOCOL NOTE

## INVITED SPECIAL ARTICLE

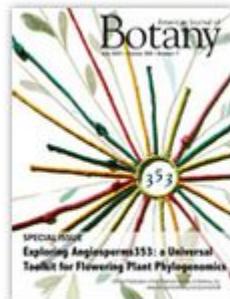
*For the Special Issue: Exploring Angiosperms353: a Universal Toolkit for Flowering Plant Phylogenomics*

## The best of both worlds: Combining lineage-specific and universal bait sets in target-enrichment hybridization reactions

Kasper P. Hendriks<sup>1,2,12</sup> , Terezie Mandáková<sup>3</sup>, Nikolai M. Hay<sup>4</sup>, Elfy Ly<sup>1</sup>, Alex Hooft van Huysduynen<sup>1</sup>, Rubin Tamrakar<sup>5</sup>, Shawn K. Thomas<sup>6</sup>, Oscar Toro-Núñez<sup>7</sup> , J. Chris Pires<sup>8</sup>, Lachezar A. Nikolov<sup>9</sup>, Marcus A. Koch<sup>9</sup> , Michael D. Windham<sup>4</sup>, Martin A. Lysak<sup>3</sup>, Félix Forest<sup>1,10</sup>, Klaus Mummenhoff<sup>2</sup>, William J. Baker<sup>10</sup> , Frederic Lens<sup>11</sup>, and C. Donovan Bailey<sup>5</sup>

# Comparação dos métodos

Target capture



American Journal of Botany: Volume  
108, Issue 7

Special Issue: Exploring  
Angiosperms353: a Universal Toolkit  
for Flowering Plant Phylogenomics

Pages: i-v, 1059-1306

July 2021

# Comparação dos métodos

## Target capture - Angiosperms353



HOME TREE OF LIFE SPECIES GENES MORE ACCESS FTI

Data release 3.0 (April 2023): 10,699 angiosperm specimens from 64 orders, 413 families, 8,336 genera and 10,377 species [View release history](#)

Please Note: Our SFTP is down temporarily, and some data may be unavailable to download. We apologise for any inconvenience. For further enquiries contact [treeoflife@kew.org](mailto:treeoflife@kew.org).

### Kew Tree of Life Explorer

The Kew Tree of Life Explorer is the gateway to Kew's research and data on the plant tree of life. We are building a comprehensive evolutionary tree of life for flowering plants and are sharing our results and data here.

[View All Species](#)

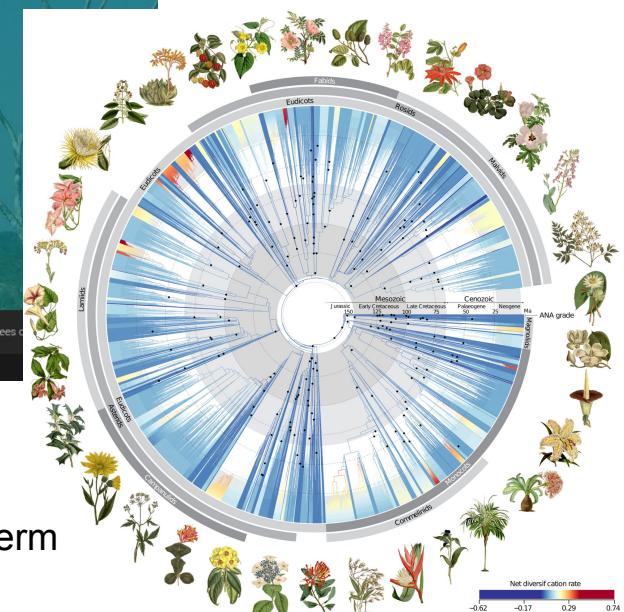
[View Tree of Life](#)

[Terms and conditions](#) [Privacy](#) [Cookies](#) [Accessibility](#) [Modern slavery](#)

Visit [kew.org](#)

<https://treeoflife.kew.org/>

8,000 (about 60%) angiosperm genera



Zuntini *et al.* Nature 629. (2024)

# Comparação dos métodos

## Target capture - dados disponíveis em bases de dados

EBI Search > Search results

EMBL's European Bioinformatics Institute

### EBI Search

Access all EMBL-EBI resources

PRJNA602196  Advanced search Examples: VAV\_HUMAN, tp53, Sulston...

Search results for **PRJNA602196**

Showing 15 results out of 626 in All results

Give us feedback on these results

Filter your results

Nucleotide sequences (626 results)

Source: ENA Study (ID: PRJNA602196)  
PRJNA602196  
Phylogenomics of Ochnaceae  
Cross references: Nucleotide sequences (625)

Source: Study (Read/Analysis) (ID: SRP243913)  
SRP243913  
Phylogenomics of Ochnaceae  
Cross references: Nucleotide sequences (626)

ENAv European Nucleotide Archive

Project: PRJNA602196  
Targeted enrichment of nuclear loci is used to resolve phylogenetic relationships of the pantropical family Ochnaceae

Secondary Study  
Accession: SRP243913  
Study Title: Phylogenomics of Ochnaceae  
Center Name: Senckenberg Research Institute and Natural History Museum  
ENA-REFSEQ: N  
PROJECT-ID: 602196  
ENA-FIRST-PUBLIC: 2021-02-03  
ENA-LAST-UPDATE: 2023-05-19

Read Files

Show Column Selection

Download report: JSON TSV Get download script Download selected files

Study Accession	Sample Accession	Experiment Accession	Run Accession	Tax Id	Scientific Name	Generated FASTQ files: FTP
PRJNA602196	SAMN13885322	SRX7618322	SRR10951301	1501052	Ochna afzelli	<input type="checkbox"/> SRR10951301_1.fastq.gz <input type="checkbox"/> SRR10951301_2.fastq.gz
PRJNA602196	SAMN13885318	SRX7618315	SRR10951308	486160	Luxemburgia bracteata	<input type="checkbox"/> SRR10951308_1.fastq.gz <input type="checkbox"/> SRR10951308_2.fastq.gz
PRJNA602196	SAMN13885317	SRX7618313	SRR10951310	2699541	Luxemburgia aff. corymbosa Cardoso	<input type="checkbox"/> SRR10951310_1.fastq.gz <input type="checkbox"/> SRR10951310_2.fastq.gz

View: XML XML (STUDY)  
Download: XML XML (STUDY)  
Navigation: Show  
Read Files: Hide  
Publications: Show  
ORCID Data Claims: Show  
Related ENA Records: Show

Tags: xref EuropePMC

# Target capture - *Angiosperms353* - Hybpiper



# **Target capture - *Angiosperms*353 - Hybpiper**

Fluxo de trabalho após sequenciamento Illumina

**Controle  
qualidade**

**Montagem  
das reads**

**Filtragem das  
sequências**

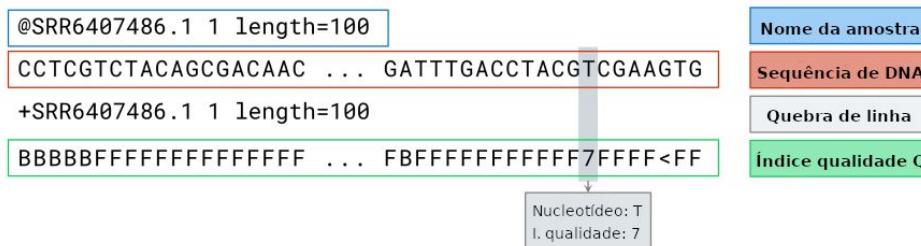
**Inferência  
filogenética**

# Target capture - Hybpiper

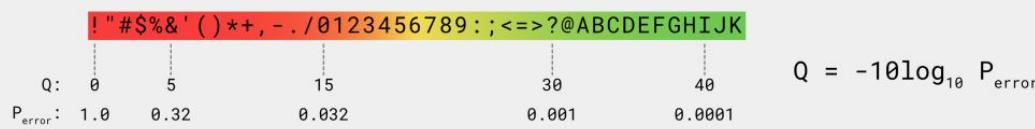
## Formato FASTQ

SRR10951301\_2.fastq.gz

## Arquivo FASTQ:



Índice de qualidade Q associado a cada um dos caracteres do ASCII



Phred Quality Score	Probability of incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%

# Target capture - Hybpiper

## Controle de qualidade das reads

SRR10951301\_2.fastq.gz



A modular tool to aggregate results from bioinformatics analyses across many samples into a single report.

### General Statistics

Showing 192/192 rows and 3/6 columns.

Sample Name	% Dups	% GC	M Seqs
SPRL_184_2	82.0%	41%	2.3
SPRL_186_1	68.9%	41%	7.0
SPRL_186_2	64.3%	41%	7.0
SPRL_187_1	80.5%	44%	6.0
SPRL_187_2	77.7%	45%	6.0
SPRL_190_1	72.5%	42%	14.9
SPRL_190_2	66.4%	42%	14.9
SPRL_192_1	76.1%	42%	11.4
SPRL_192_2	71.2%	42%	11.4

# Target capture - Hybpiper

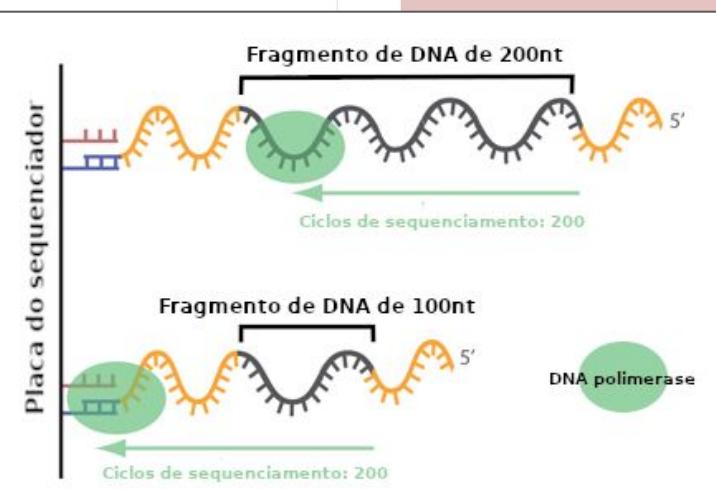
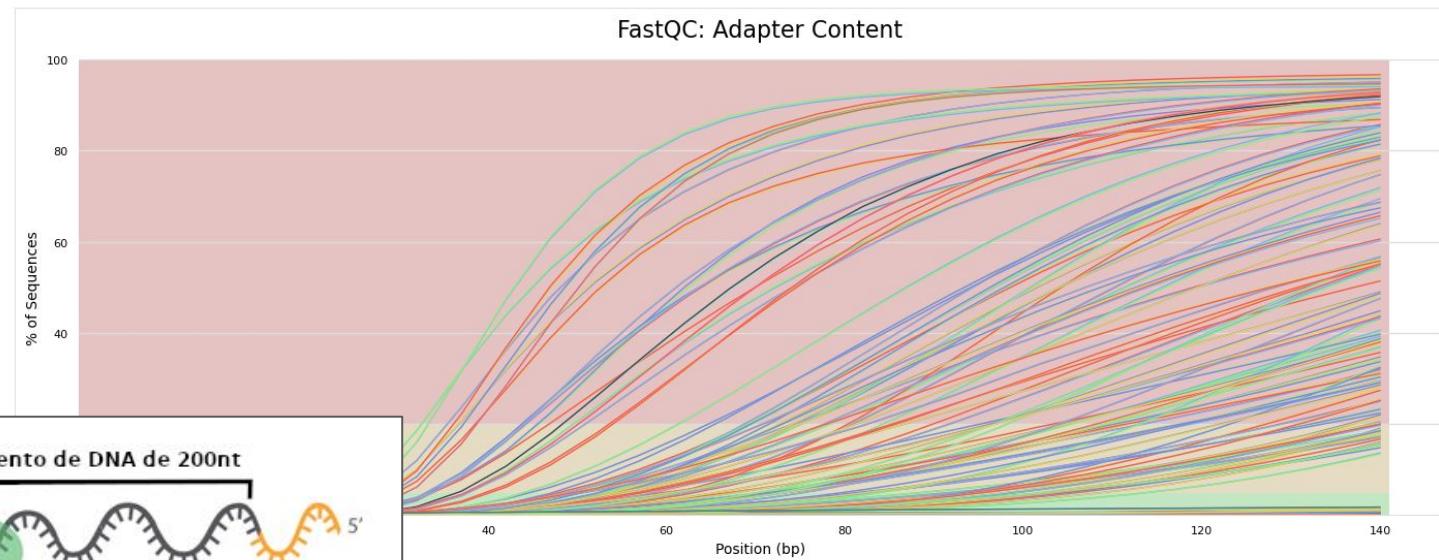
## Controle de qualidade das reads

### Adapter Content

192

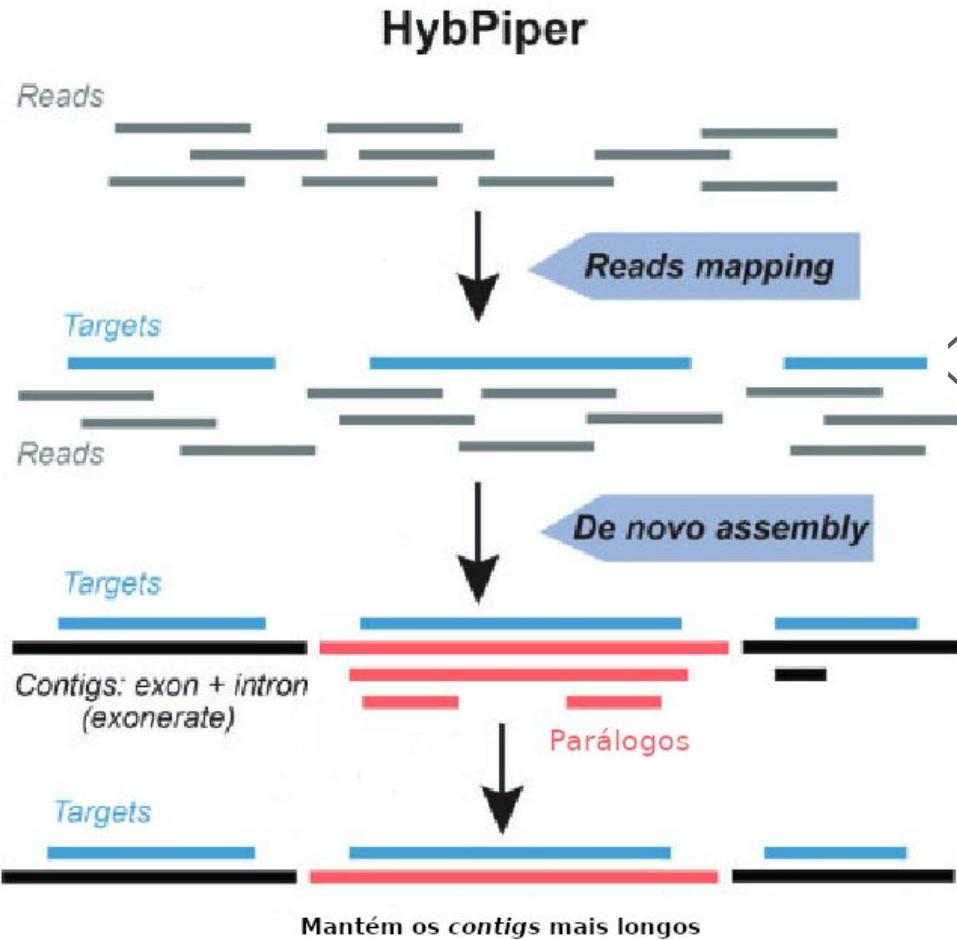
The cumulative percentage count of the proportion of your library which has seen each of the adapter sequences at each position.

Flat image plot. Toolbox functions such as highlighting / hiding samples will not work (see the docs).



# Target capture - Hybpiper

## Montagem das sequencias



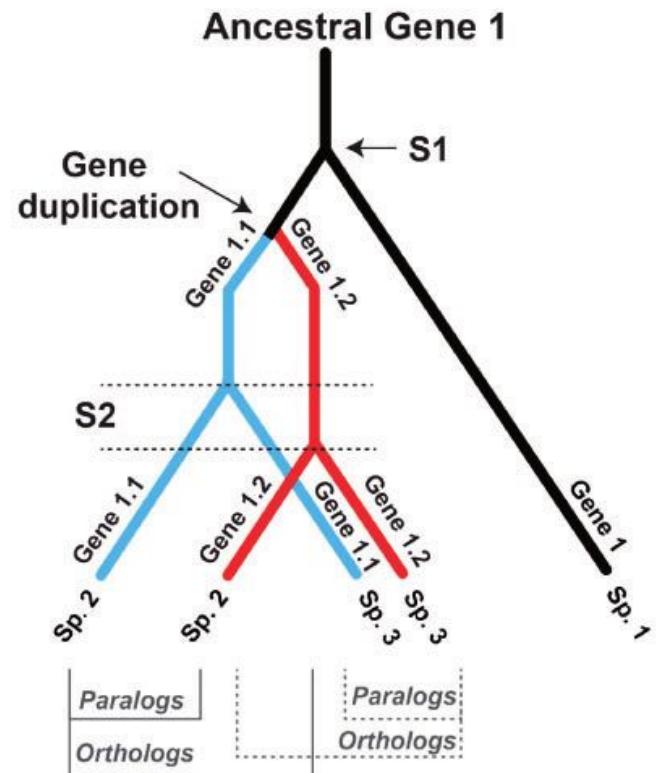
Consensus:	g a g g a t g a t g a t g t t g t t g t t g t t g
AJFN-4471	[ 1 A A T G T T A T C A C A G G A T G A A G A G A G A
Ambrt-4471	[ 1 A G T G T T A T T C A A G A T G A A G A T G
Arath-4471	[ 1 A A C G T G G T T G A A G A T G A A G A A A
CCID-4471	[ 1 A A T G T C C T T C A C G G A T G A A G A T A
IDAU-4471	[ 1 A A T G T A A T C C A C G G A T G A A G A T A
NHUA-4471	[ 1 A A T G T G A T T G A A G A T G A A G A C G
NNOK-4471	[ 1 G G T G T A G T T C A G G A C G A A G A C A
Orysa-4471	[ 1 G G T G T G G T T C A C G G A T G A A G A C A
QUTB-4471	[ 1 G T T T C T T C T G G A T A A T T C A A G
TJLC-4471	[ 1 G T A T G T G A A G C A G G A T G C A T A T
TVSH-4471	[ 1 G T T C T C A G A T T C T T A A A T T
VUSY-4471	[ 1 A A G T T A A T C C A G G A T G A A G A G A
Ambrt-4527	[ 1 G A G G A G C G G G T G A T T G C C T T G
Arath-4527	[ 1 G A A G A G A G A G T G A A T G T T C T T G
BERS-4527	[ 1 G C T C C A C T T G T T G C T G C T C T
GRFT-4527	[ 1 G G A A G G G A G C A T T G C C T T G T T G
HAEU-4527	[ 1 G A A G A G A G G G T G G T T G T G T T G
JEPE-4527	[ 1 G A A G A A A G A G T T T C T G T A T T G G
NUZN-4527	[ 1 G A A G G T T C C A A G A A C T T C A T G A
RCUX-4527	[ 1 G A A G A G A G G G T T G T G G T G T T G
TIUZ-4527	[ 1 G A A G A G A G G G T T A A T T G T G C T A G
ZENX-4527	[ 1 G A T T G A G A G G A T A A C T G T A T T G G
Ambrt-4691	[ 1 C A A C T T G A A G A A T G T A G C C T G T A
Arath-4691	[ 1 G A G C T T G G A T C T A T C G C C T T G C G
BHYC-4691	[ 1 G A A G T T G A A C T T G T G T T G C T T G T G
IDAU-4691	[ 1 G A G A T T C C G G A G C T A G C C T T G C G
KXSK-4691	[ 1 G G A A T T G A G A A A T G T T G C T T G T G
Orysa-4691	[ 1 G C G G T T G G G C G G G C T T G C G T G C G

Target file

# Target capture - Hybpiper

Filtragem das sequências - Parálogos

- Alelos
- Verdadeiros parálogos
- Contaminação



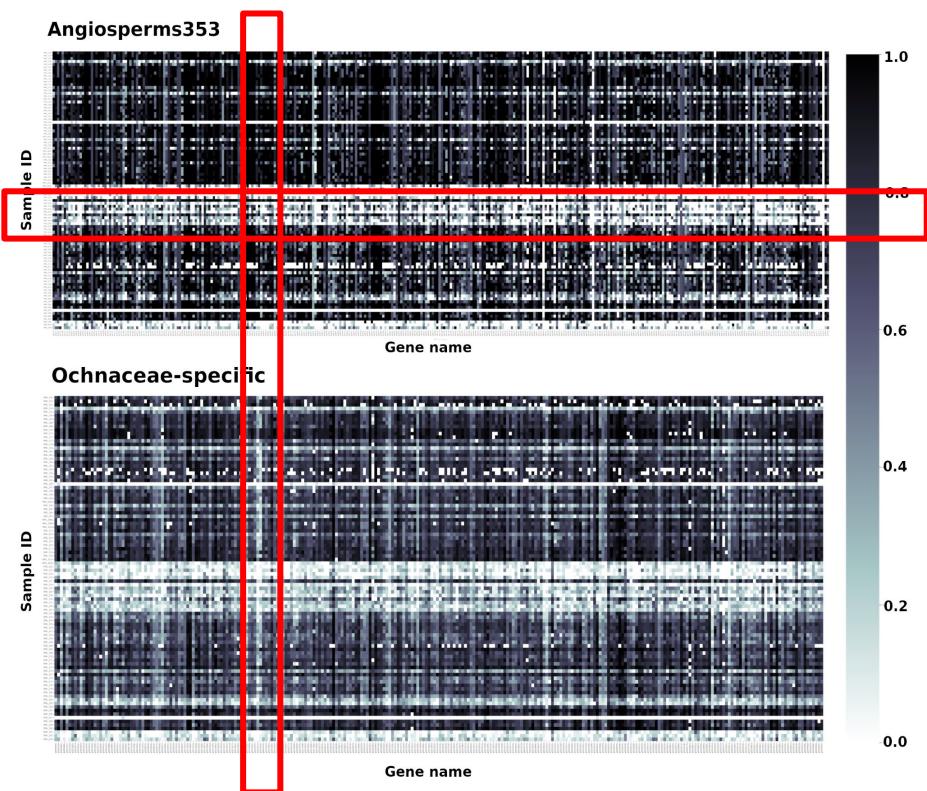
Zhou *et al.* *Sys. Bio.* 71: 410–425 (2022)

# Target capture - Hybpiper

Filtragem das sequências - Dados faltantes

Baixa recuperação:

- Regiões/loci com poucas amostras
- Amostras com sequências para poucos loci
- Sequências muito curtas



# Inferência filogenética

Material para outro curso

- Sumarização de árvores de genes
- Concatenação
- Redes filogenéticas

## Accurate Species Tree EstimatoR (ASTER\*)



A family of optimization algorithms for species tree inference:

1. [ASTRAL-IV](#) (from unrooted gene tree **topologies** with integrated CASTLES-II for branch lengths)
2. [ASTRAL-Pro3](#) (from unrooted gene **family** tree topologies with integrated CASTLES-Pro)
3. [Weighted ASTRAL](#) (from unrooted gene trees with branch **lengths** and/or **supports**)
4. [CASTER-site](#) (from **whole genome alignments** or aligned sequences)
5. [CASTER-pair](#) (from **whole genome alignments** or aligned sequences)
6. [WASTER-site](#) (from **raw reads**)
7. SISTER (from optical-map-like distance data or shape data)
8. MONSTER



# Considerações finais

- NGS ou HTS tem o melhor custo benefício para estudos em sistemática hoje.
- Target sequence é uma ótima opção para inferência filogenética a diferentes níveis, funciona muito bem em material de herbário (e.g. tipos).
- Mais dados não representam maior resolução mas sim a possibilidade de explorar fenômenos de evolução molecular interessantes (e.g. hibridização, duplicação genica).
- Estratégias de limpeza e filtragem dos dados são fundamentais para uma ‘boa’ reconstrução filogenética.

# Para saber mais:



Curso completo de Introdução a Next Generation Sequencing (NGS).

Profa. Dr. Carolina M. Siniscalchi  
Mississippi State University



Slides, lista com materiais de apoio:

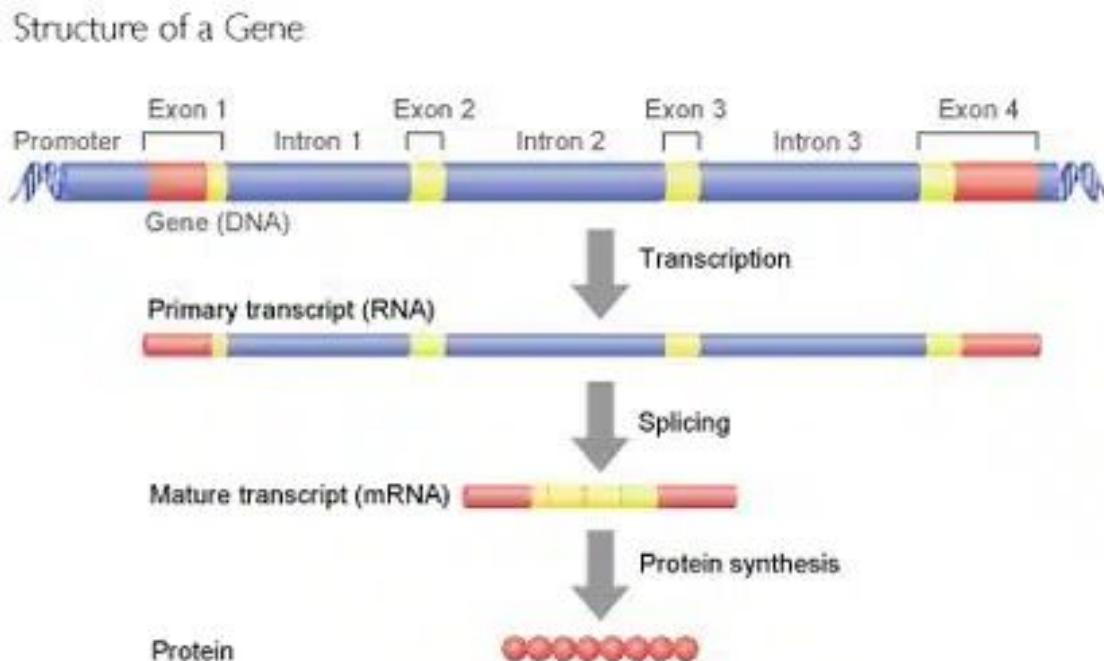
- Bash, R e Python para iniciantes
- Tutoriais
- Livros



**Dúvidas?**

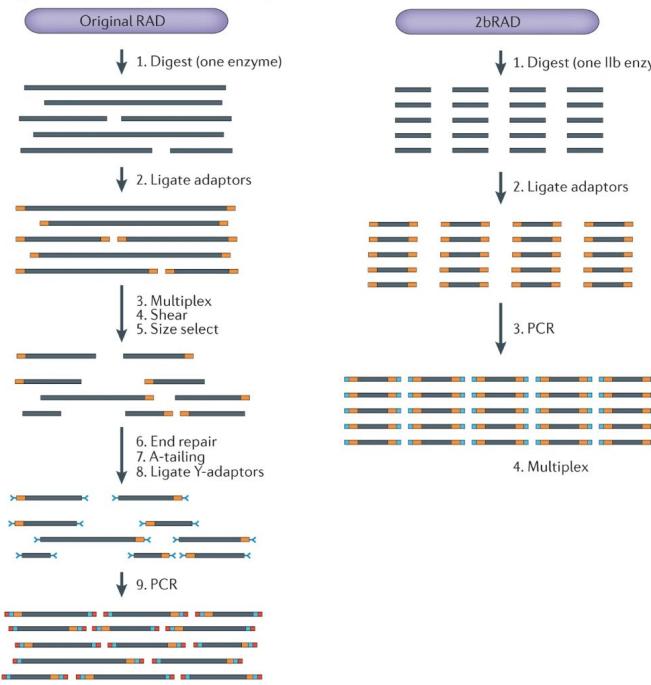
# Target capture - Hybpipe

## Montagem das sequencias

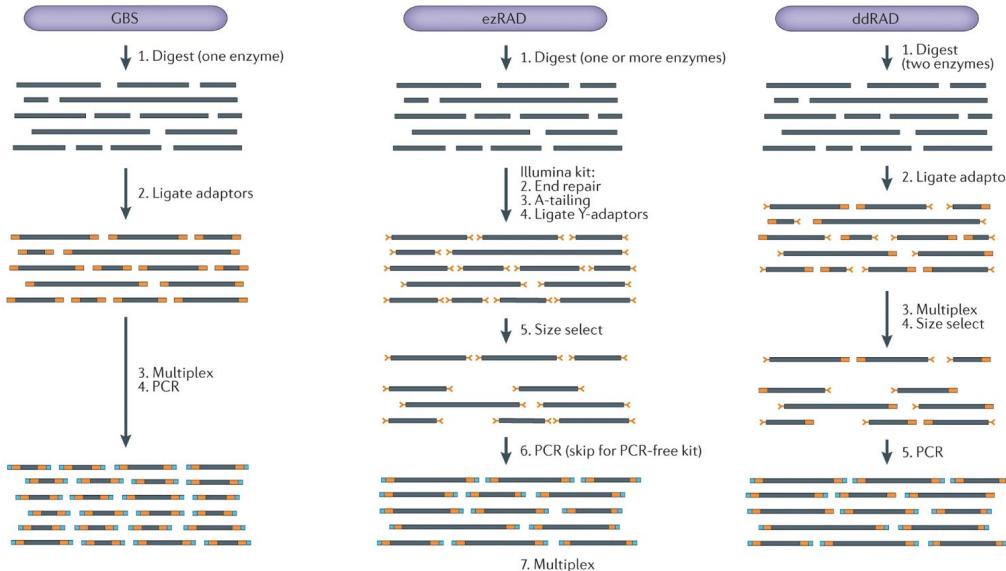


© Wellcome Trust

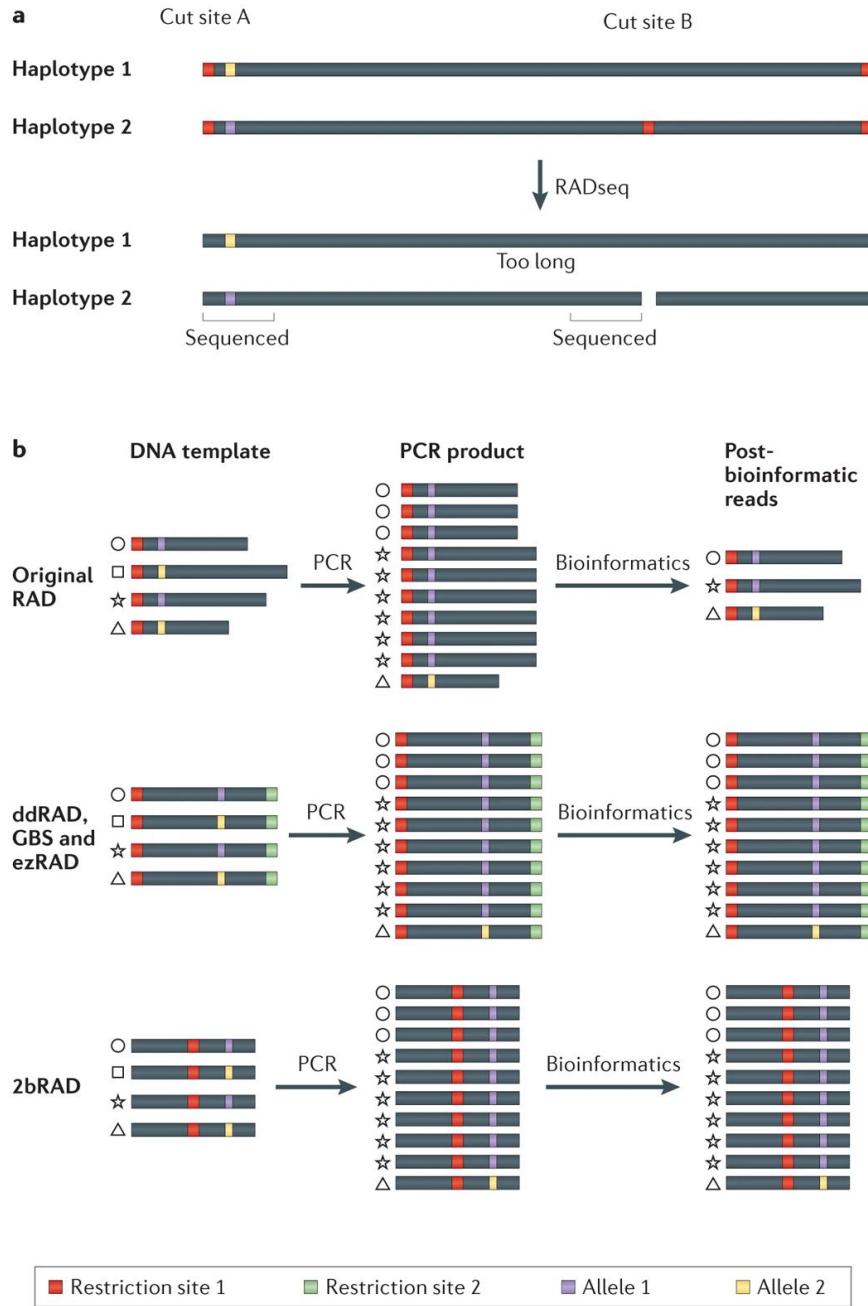
### Sequence next to single restriction enzyme cut sites



### Sequence flanked by two restriction enzyme cut sites



All protocols begin by digesting high-molecular-weight genomic DNA with one or more restriction enzymes. For most protocols, the sequencing adaptors (oligonucleotides) are added in two stages, with one set of oligonucleotides added during a ligation step early in the protocol, and a second set of oligonucleotides incorporated during a final PCR step. The second set of oligonucleotides extends the length of the total fragment to produce the entire Illumina adaptor sequences. By contrast, the original RADseq adds adaptors in three stages. For Illumina sequencing, the adaptors on either end of each DNA fragment must differ, and therefore some protocols (for example, original RADseq, double digest RAD (ddRAD) and ezRAD) use Y-adaptors that are structured to ensure that only fragments with different adaptors at either end are PCR-amplified (illustrated here as Y-shaped adaptors). Other protocols (for example, genotyping by sequencing (GBS)) simply rely on the fact that fragments without the correct adaptors will not be sequenced. To generate fragments of an ideal length for sequencing, most methods use common-cutter enzymes (for example, 4–6 bp cutters) to generate a wide range of fragment sizes, followed by a direct size selection (gel-cutting or magnetic beads, for example, ezRAD and ddRAD) or an indirect size selection (as a consequence of PCR amplification or sequencing efficiency, for example, GBS).



An example of allele dropout for a restriction site-associated DNA sequencing (RADseq) protocol that uses size selection to reduce the number of loci to be sequenced. Grey lines represent chromosomes within one individual, red squares represent restriction cut sites, coloured squares represent heterozygous SNPs, and square brackets represent genomic regions that are sequenced. Mutation in cut site B for haplotype 1 makes the post-digestion fragment containing the SNP too long to be retained during size selection for haplotype 1, eliminating the possibility of sequencing of any loci on that fragment, and causing the individual to appear homozygous at the heterozygous SNP. **b** | An example of fragments produced after PCR for one heterozygous locus for different RADseq protocols, and the reads retained after bioinformatic analyses. PCR duplicates are shown with the same symbol (circle, square, asterisk or triangle) as the parent fragment from the original template DNA. By chance, some alleles will amplify more than others during PCR. For all protocols, PCR duplicates will be identical in sequence composition and length to the original template molecule. For the original RADseq, this feature (that is, identical length) can be used to identify and remove PCR duplicates bioinformatically, because original template molecules for a given locus will not be identical in length. For alternative RADseq methods, this feature cannot be used to identify PCR duplicates, because all original template molecules for a given locus are identical in length. High frequencies of PCR duplicates can cause heterozygotes to appear as homozygotes or can cause PCR errors to appear as true diversity. Part **b** is adapted with permission from Ref. 37, Wiley.

**A**

<i>SPRL_187 multi_hit_stitched_contig_comprising_3</i>	49	G G T A A T A C A T T G G G T C C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_168.main NODE_1_length_2704_cov_301</i>	49	G G C A A T A C A T T G G G T A C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_168.0 NODE_2_length_2204_cov_71.248</i>	1	G G T A A T A C A T T G G G T C C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_190.0 NODE_2_length_2408_cov_449.02</i>	49	G G T A A T A C A T T G G G T C C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_190.main NODE_1_length_2824_cov_540</i>	49	G G C A A T A C A T T G G G T A C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_201ii.main NODE_1_length_2542_cov_18</i>	10	G G C A A T A C A T T G G G T A C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A
<i>SPRL_201ii.0 NODE_2_length_2054_cov_49.41</i>	1	G G T A A T A C A T T G G G T C C C G A G G T C A A T A C C A T C A G C T T G A A G A G G A A G A

**B**

# What is long-read sequencing?

Long-read sequencing is a DNA sequencing technique that enables the sequencing of much longer DNA fragments than traditional short-read sequencing methods. While short reads can capture the majority of genetic variation, long-read sequencing allows the detection of complex structural variants that may be difficult to detect with short reads. These include large inversions, deletions, or translocations, some of which have been implicated in areas like genetic disease.

Long-read technology can help resolve challenging regions of the genome by sequencing thousands of bases to:

- Resolve traditionally difficult to map genes or regions of the genome, such as those containing high variable or highly repetitive elements
- Perform phased sequencing to identify co-inherited alleles, haplotype information, and phase *de novo* mutations
- Generate long reads for *de novo* assembly and genome finishing applications

# Técnicas de sequenciamento

Sequenciamento de alto rendimento - HTS (2005 - presente)

- Sequencia milhões de fragmentos simultaneamente (em paralelo) em cada corrida (Número dependendo da plataforma)
- Custo por nucleotídeo sequenciado é muito menor quando comparado com Sanger
- Precisa de preparação de bibliotecas genômicas

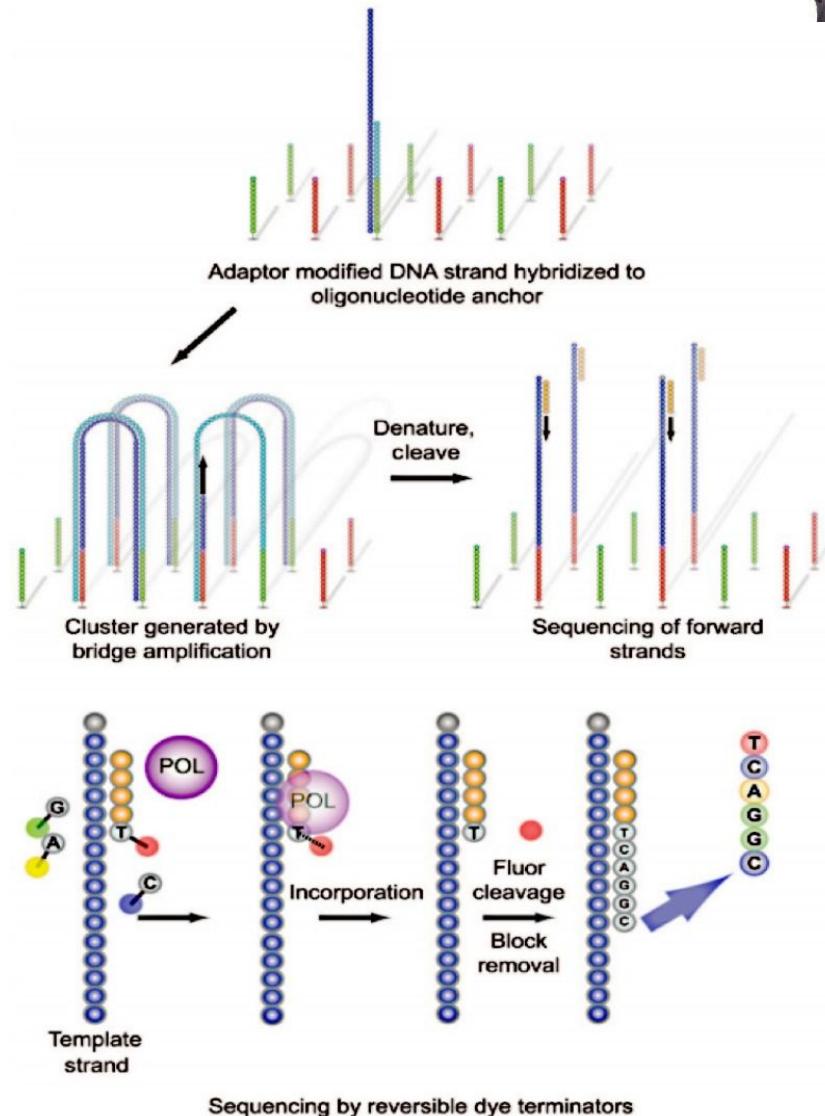
# Técnicas de sequenciamento

## Sequenciamento Illumina



### Três etapas principais:

1. Ligação das bibliotecas na plataforma
2. PCR fase sólida (Bridge amplification)
3. Sequenciamento por síntese



<https://www.youtube.com/watch?v=fCd6B5HRaZ8>