

Genomics on a Shoestring Budget

Thomas Gatter

17.02.2023

Bioinformatics, Leipzig University



Can we study the genome of novel organisms on a small budget?

Can we study the genome of novel organisms on a small budget?

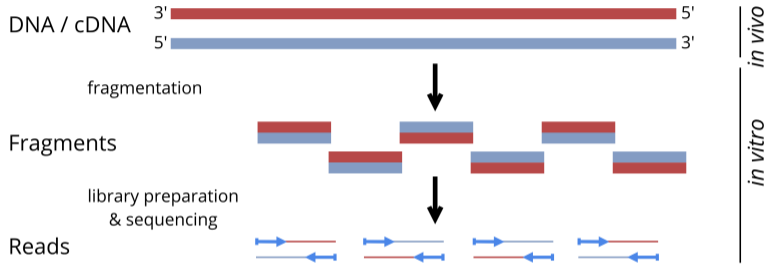
Don't do it. You just get a lot of money or it won't be worth it.
But that's also way too much spending for this proposal. So you
won't get it here.

– Every Reviewer

Proposed Solution:

Get a lot of Money (somehow) → Sequence with PacBio HiFi + Hi-C
→ Process with Super-Computer → Profit (?)

Option 1: Short Read Sequencing

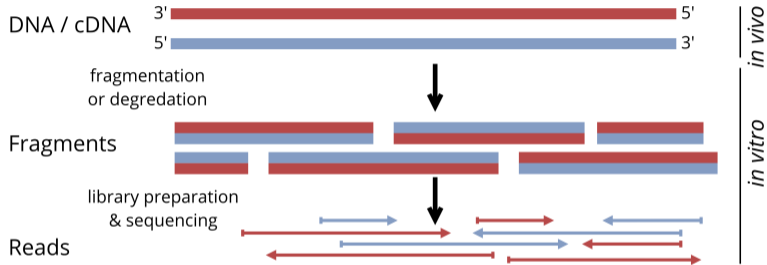


Read length: $2 \times 75-150$ bp

Accuracy: $> 99\%$

Provider: Illumina, ThermoFischer, Roche

Option 2: Long Read Sequencing



Read length: avrg. > 10 kbp, often much longer

Accuracy: $\approx 90\%$ \rightarrow better with newer iterations

Provider: Oxford Nanopore(, Pacific Biosciences)

Short Read Assembly

- + cheap
- + accurate
- + well established with de Bruijn or String-Graphs
- cannot resolve longer repeats

Short Read vs Long Read Assembly

Short Read Assembly

- + cheap
- + accurate
- + well established with de Bruijn or String-Graphs
- cannot resolve longer repeats

Long Read Assembly

- + resolve long complex regions
- less accurate
- more expensive
- expensive pairwise overlapping
- require high coverage

Short Read vs Long Read Assembly

Short Read Assembly

- + cheap
- + accurate
- + well established with de Bruijn or String-Graphs
- cannot resolve longer repeats

Long Read Assembly

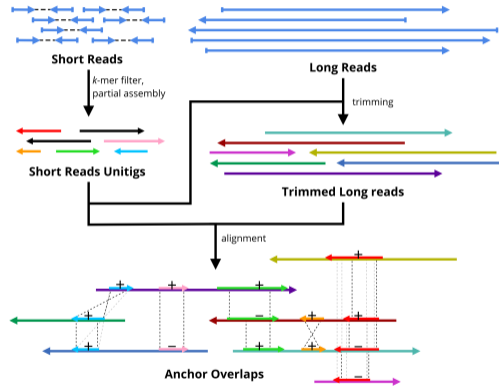
- + resolve long complex regions
- less accurate
- more expensive
- expensive pairwise overlapping
- require high coverage

Combine both!

The LazyB Workflow



Short reads are assembled to build accurate anchors between long reads.

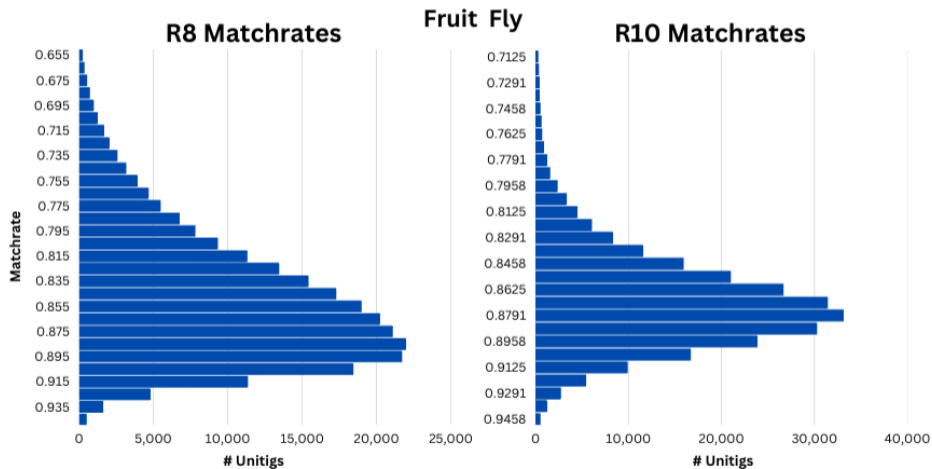


Then some ma(th)gic happens...

Fruit Fly

X	Tool	compl.[%]	#ctg	#MA	NA50
~5x	LazyB	71.624	1879	68	64415
	Canu	-	-	-	-
	Wtdbg2	6.351	2293	2	-
	HASLR	24.484	1407	10	-
	DBG2OLC	25.262	974	141	-
	Wengan	81.02	2129	192	77215
~10x	LazyB	80.111	596	99	454664
	Canu	49.262	1411	275	-
	Wtdbg2	41.82	1277	155	-
	HASLR	67.059	2463	45	36979
	DBG2OLC	82.52	487	468	498732
	Wengan	84.129	926	237	221730
~45x	ABYSS	83.628	5811	123	67970

Advances in Nanopore Sequencing



Fruit Fly - R 8 vs R10

X	Tool	Chem.	compl.[%]	#ctg	#MA	NA50
~5x	LazyB	R8	71.624	1879	68	64415
		R10	71.028	708	91	189244
	DBG20LC	R8	25.262	974	141	-
		R10	33.413	895	161	-
~10x	Wengan	R8	81.02	2129	192	77215
		R10	78.564	1645	140	117504
	LazyB	R8	80.111	596	99	454664
		R10	78.206	191	91	1031893
	DBG20LC	R8	82.52	487	468	498732
		R10	87.519	230	281	1016141
	Wengan	R8	84.129	926	237	221730
		R10	83.037	483	182	528879
~45x	ABySS		83.628	5811	123	67970

Meccus longipennis



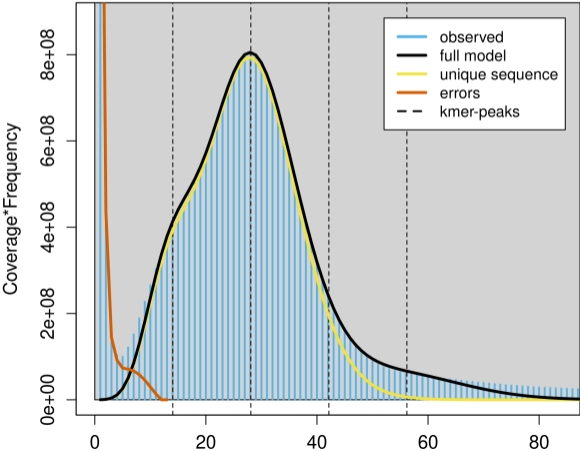
- blood sucking bug
- native to latin america
- host for Trypanosoma cruzi
- vector for Chagas disease

Genome virtually unexplored...

Lets get practical...

GenomeScope Profile

len:1,094,642,408bp uniq:60.8%
aa:98.7% ab:1.31%
kcov:14 err:0.437% dup:0.938 k:21 p:2



Estimated genome size of 1.1 Gb

⇒

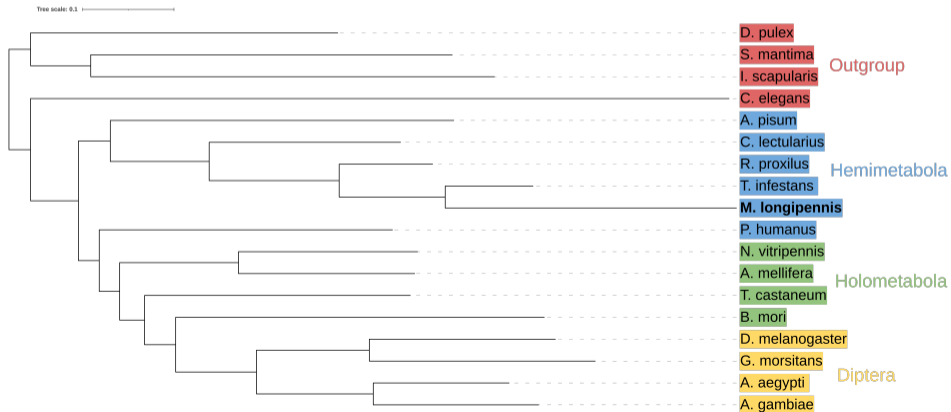
6.2× Nanopore
40× Illumina

Assembly Statistics for *M. longipennis* (with BUSCO Hemiptera ODB10)

Tool	#ctg	Assembled bp	N50	BUSCO C + F
DBG20LC	-	-	-	-
Wengan	677	5,681,409	8,642	50 (2.0%) + 2 (0.1%)
HASLR	68,585	416,784,090	8,326	1280 (51.0%) + 232 (9.2%)
ABySS	695,368	893,209,008	1,582	1314 (52.4%) + 500 (19.9%)
LazyB	48,074	788,046,408	22,713	1596 (63.5%) + 181 (7.2%)

Annotation

- RepeatModeler + BRAKER2 → 19353 Proteins
- OrthoFinder for related organisms → 7592 Orthologues to *R. prolixus*
- including orthologues for hematophagy und immune related proteins



Thank You!

Peter Stadler

Javier T Granados Riverón

Sarah von Löhneysen

Kevin Klein

Felix Kühnl

DFG
Deutsche
Forschungsgemeinschaft

de  **NBI**
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

SAB 
SÄCHSISCHE
AUFBAUBANK