

BIOINFORMATICS
INSTITUTE

Расширение функциональности инструмента для построения функции выравниваемости генома

Руководитель: Евгений Бакин
(Институт биоинформатики)

Студент: Елизавета Скалон

Genome Mappability Score (GMS)

отражает вероятность того, что рид может быть однозначно выровнен на данную позицию генома

```
Reference ...TCCTAATCGTATCTAGGCTCGATTCCGTATGATTCCGAAACG ... AACGTCTCTGTTAGGTTCTCGTATCTAGGCTCGTATAGCTAGCGTCA...
          TCGTATCTAGGCTCGATTCCGTA
          CGTATCTAGGCTCGATTCCGTAT
          GTATCTAGGCTCGATTCCGTATG
          TATCTAGGCTCGATTCCGTATGA
          ATCTAGGCTCGATTCCGTATGAT
          TCTAGGCTCGATTCCGTATGATT
          CTAGGCTCGATTCCGTAAGATTCC
          TAGGCTCGATTCCGTATGATTCC
          AGGCTCGATTCCGTATGATTCCG

          ...
          TTAGGTTCTCGTATCTAGGCTCG
          TAGGTTCTCCTATCTAGGCTCGT
          AGGTTCTCCTATCTAGGCTCGTA
          GGTTCCGATCTAGGCTCGTAT
          GTTCCGATCTAGGCTCGTATA
          TTCCCTATCTAGGCTCGTATAG
          TCCTATCTAGGCTCGTATAGC
          CTCCTATCTAGGCTCGTATAGCT
          TCGTATCTAGGCTCGATTCCGTA
```

GMS **низкий** — на данную позицию генома может выровняться много ридов
→ на данную позицию может быть выровнен неправильный рид
→ **неверная** интерпретация результатов выравнивания

GMS **высокий** — данная позиция генома уникальна
→ выравниванию **можно доверять**

Существующие реализации:

Genome Mappability Analyzer

[Bioinformatics](#). 2012 Aug 15; 28(16): 2097–2105.

PMCID: PMC3413383

Published online 2012 Jul 4. doi: [10.1093/bioinformatics/bts330](https://doi.org/10.1093/bioinformatics/bts330)

Genomic dark matter: the reliability of short read mapping illustrated by the genome mappability score

[Hayan Lee](#)^{1,2,*} and [Michael C. Schatz](#)^{1,2}

«GMStool»

Hardware-effective Tool for Genome Mappability Score Estimation

Feb 2017

Jun 2017

Евгений
Бакин
Наталья
Зорина

Александр
Предеус

Цель:

Расширить функциональность тула в части работы с различными форматами данных

Задачи:

- Добавить возможность выводить данные в форматы BED, bigBed и TDF
- Ускорить процесс выгрузки данных
- Реализовать запись данных непосредственно в формат bigWig

ИТОГИ

Разнообразие форматов

Раньше:

- Wig
- bigWig

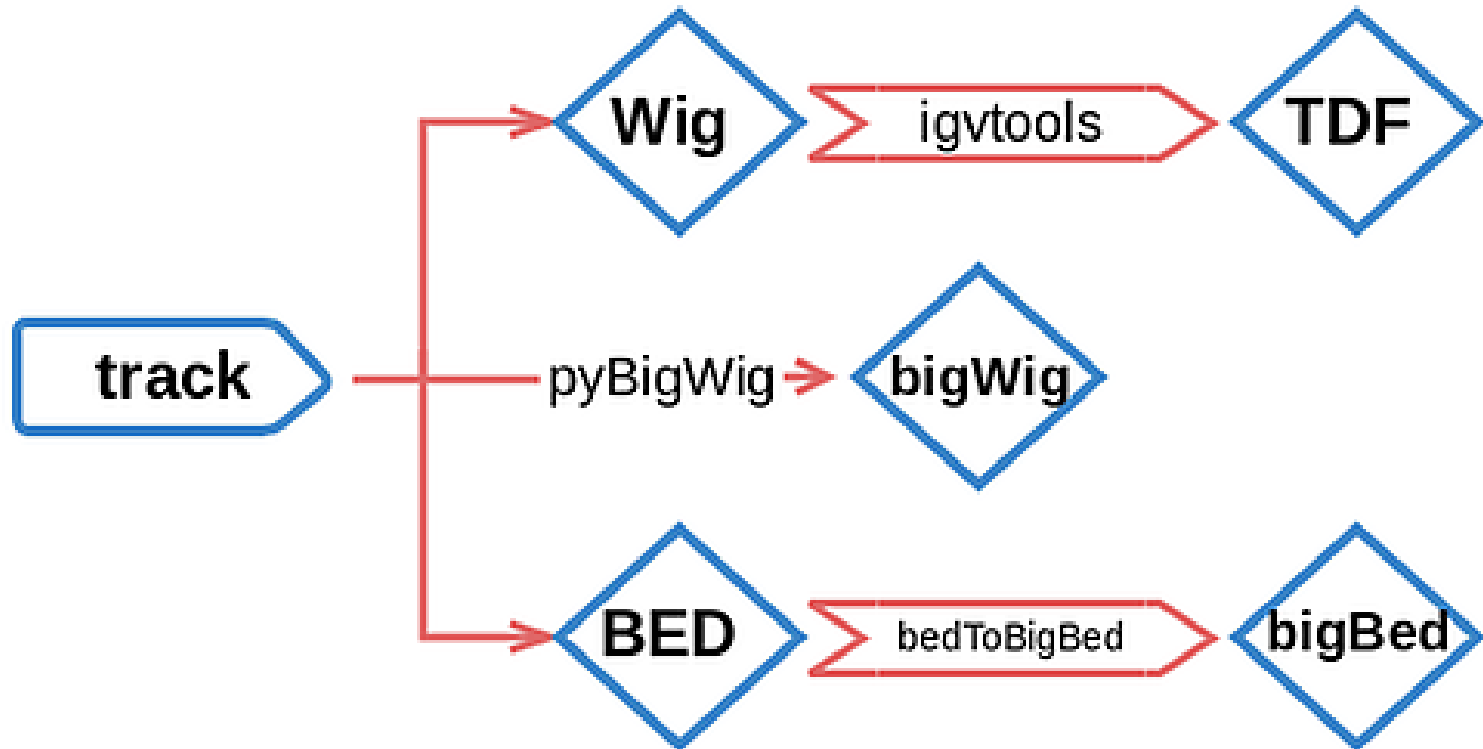


ИТОГИ

Разнообразие форматов

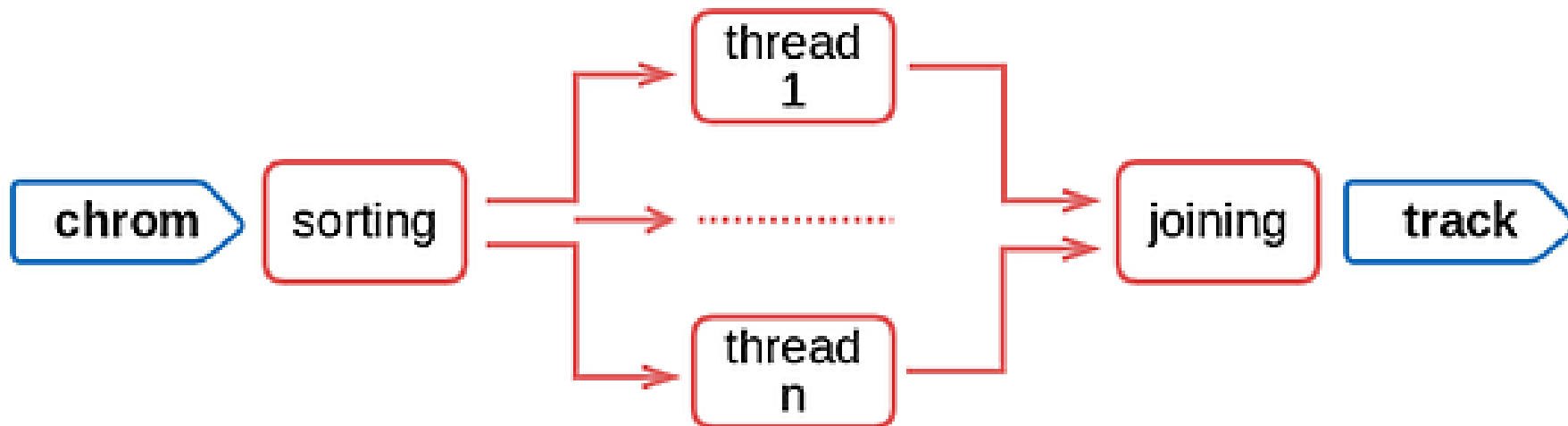
Сейчас:

- Wig
- bigWig
- BED
- bigBed
- TDF



ИТОГИ

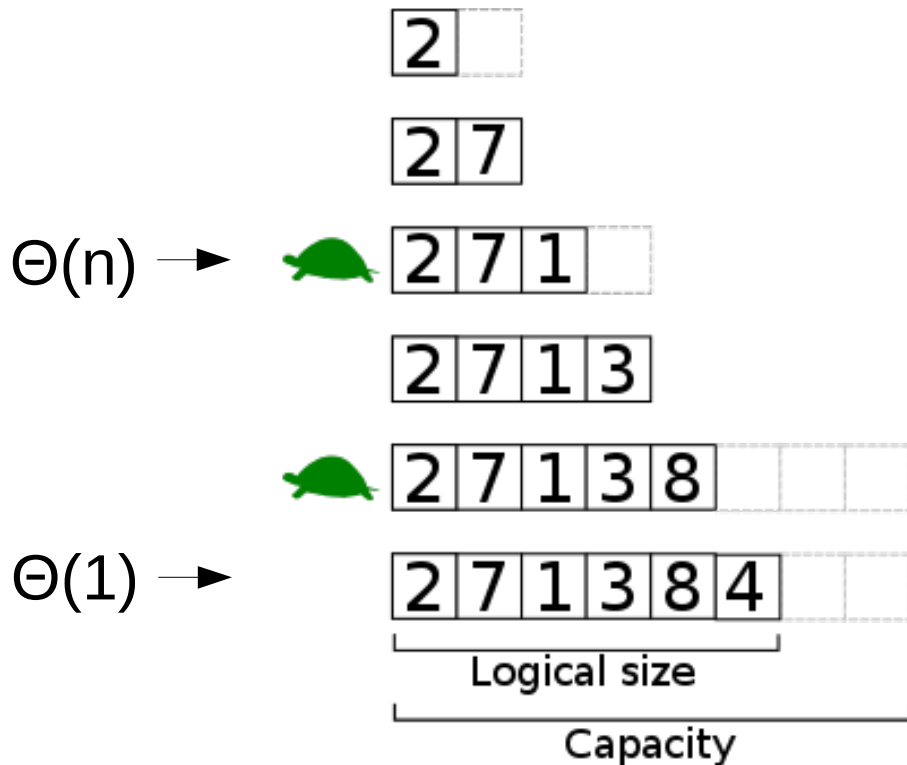
Скорость: multiprocessing



* n, то есть количество потоков, задается пользователем.

Итоги

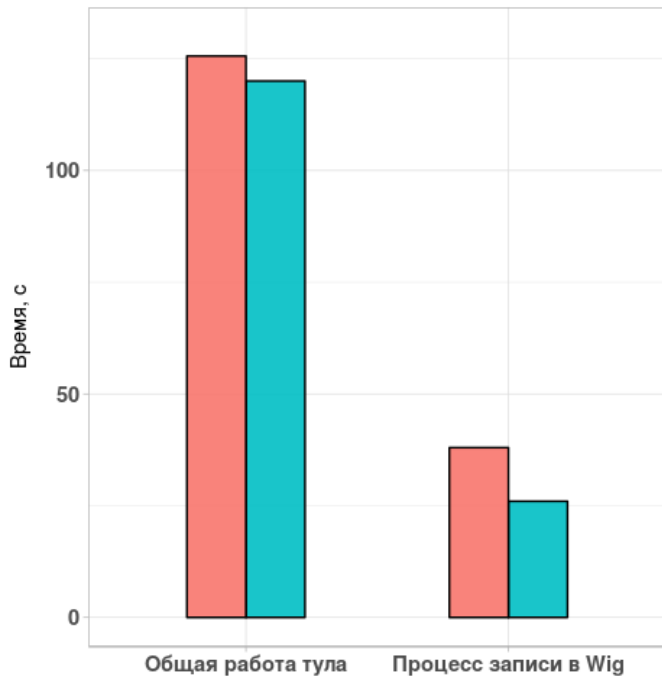
Скорость: геометрическое расширение памяти



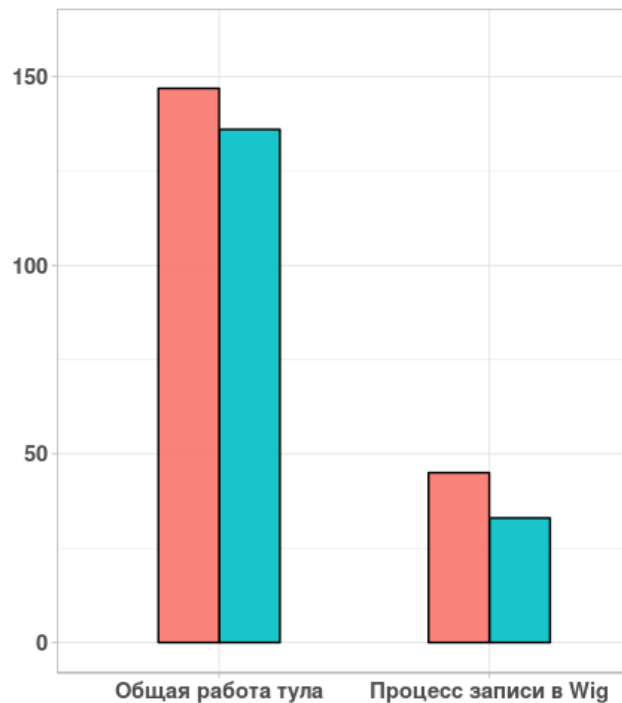
ИТОГИ

Увеличение скорости на примере экспорта в Wig

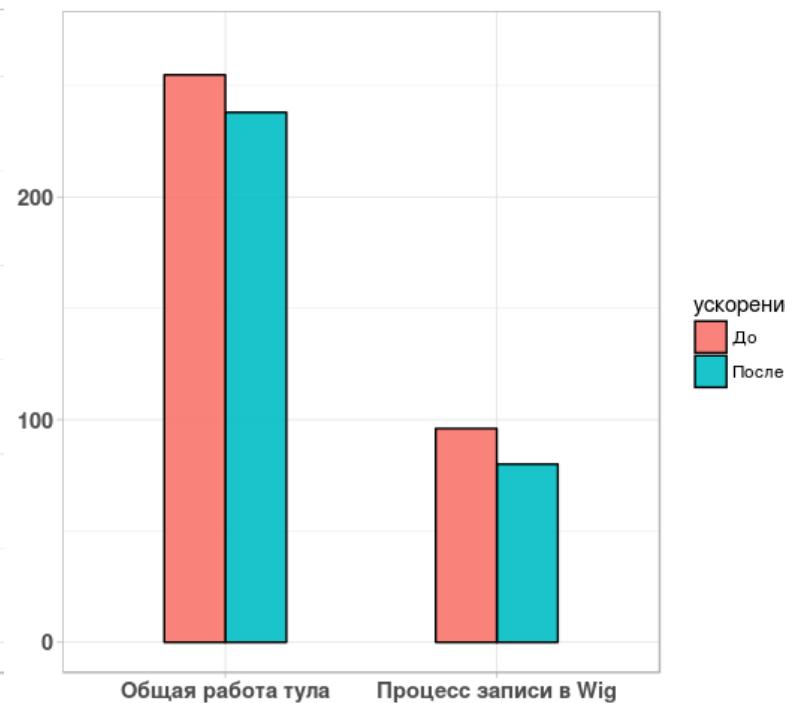
C.elegans



A.thaliana



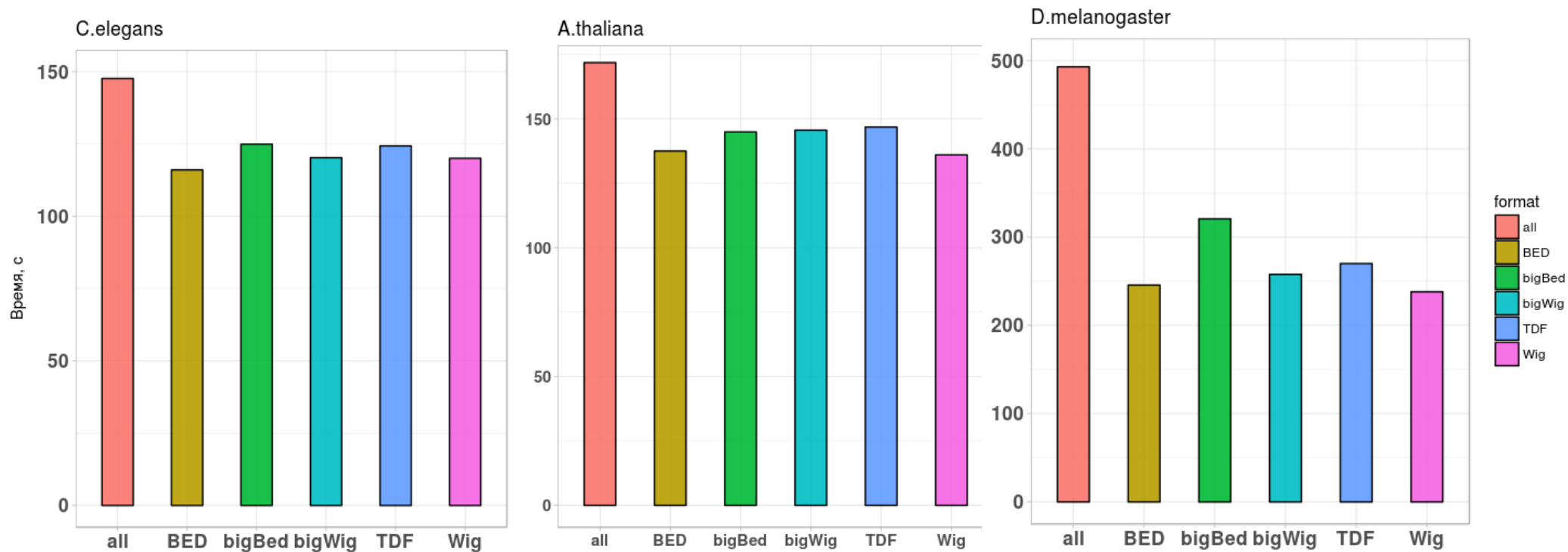
D.melanogaster



* подсчет проводился на 4-ядерном процессоре

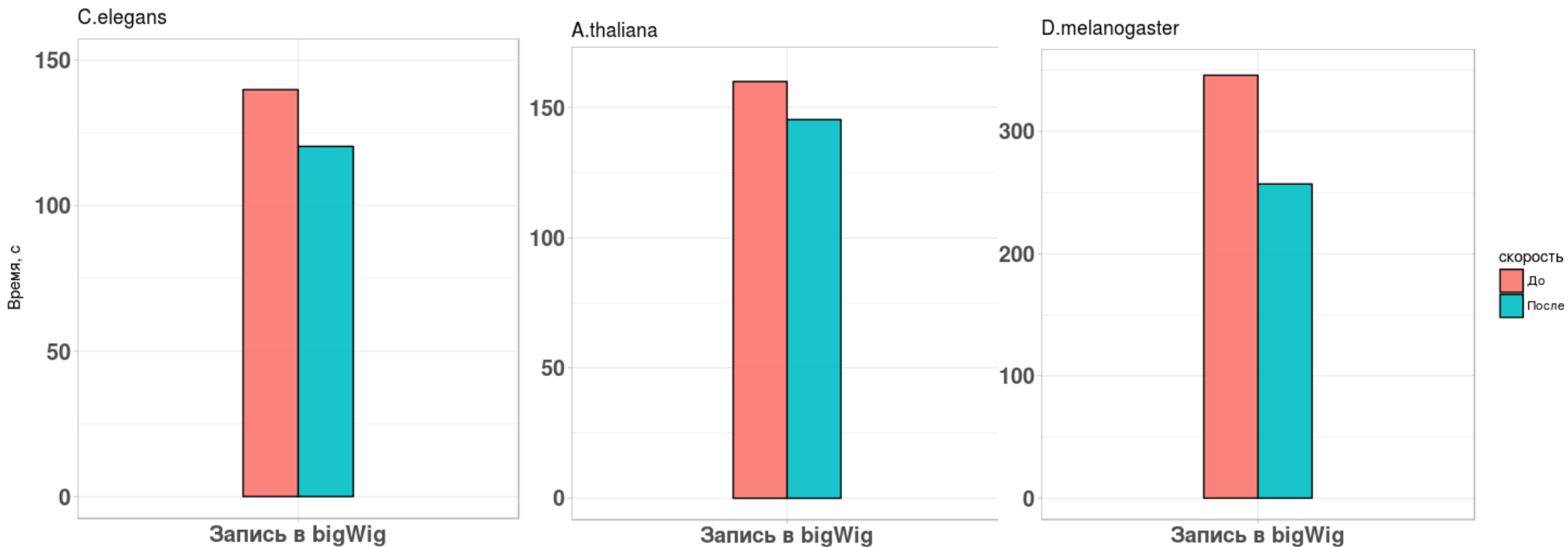
ИТОГИ

Скорость: формирование единого трека для последующей записи в разные форматы

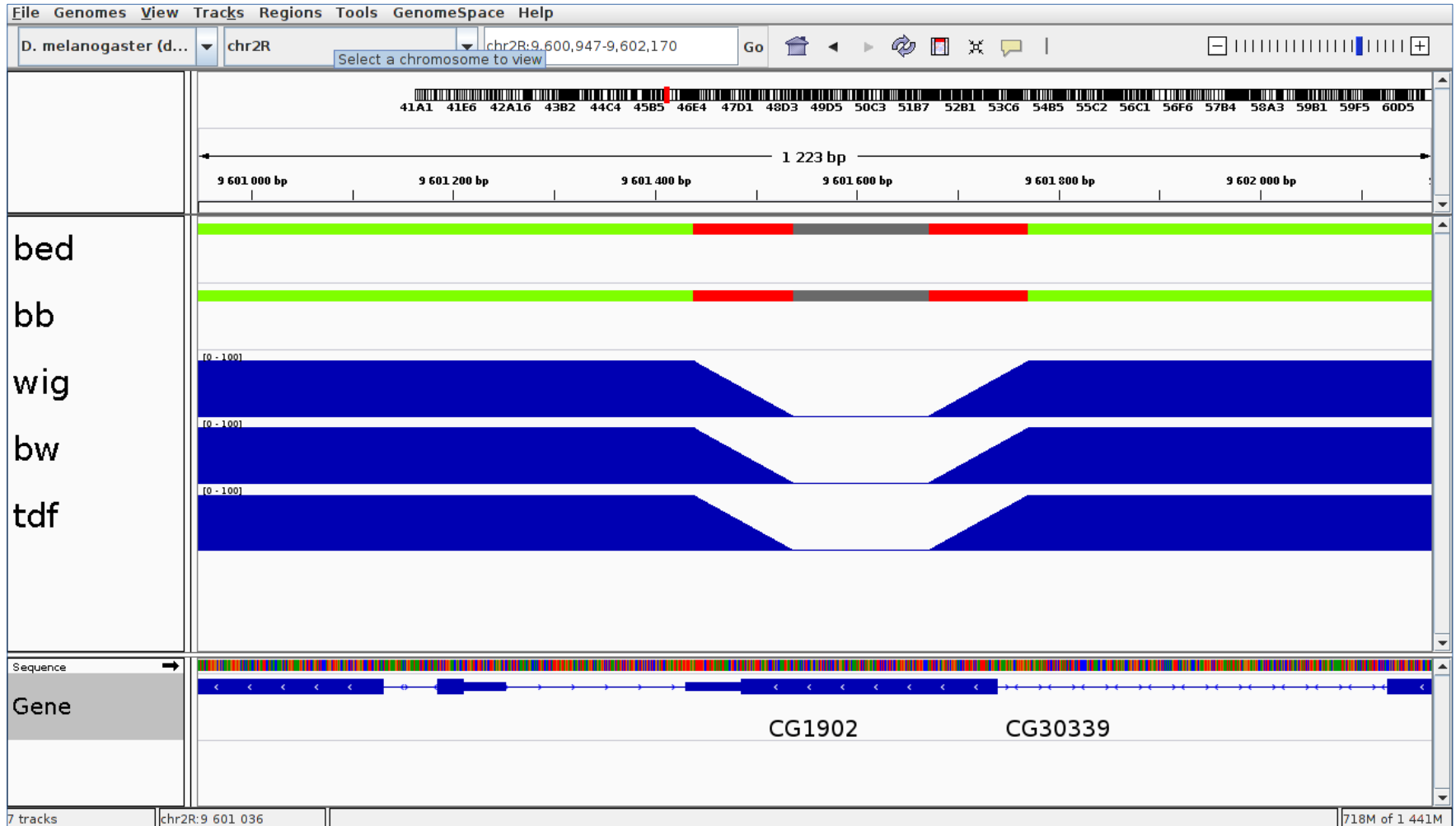


ИТОГИ

Скорость: экспорт напрямую в bigWig



Пример работы тула



Выводы

Расширен функционал тула в части экспорта полученных результатов:

- Сейчас возможен вывод в 5 различных форматах: Wig, bigWig, BED, bigBED, TDF. Каждый из них имеет свои преимущества для пользователя.
- Скорость выгрузки уменьшена за счет реализации многопоточного режима, а также за счет геометрического расширения памяти.
- Реализован экспорт данных непосредственно в формат bigWig

Планы на будущее:

- Усовершенствовать остальную часть тула (мультипроцессинг и тд.)

Спасибо за внимание!