

# Обзор литературы по ассемблерам

Дмитрий Яковлев

EPAM Systems

17 октября 2016 г.

- 1 Введение
- 2 Модель ошибки на данных BioNano
- 3 Ассемблеры
  - TWIN
  - OPTIMA
  - MAligner
  - OMBlast

# Введение



# Модель ошибки на данных BioNano



# Ассемблеры

TWIN





# OPTIMA



# MAligner



# OMBlast



Этапы выравнивания:

- Поиск стартовых мест (сидов) для начала выравнивания
- Расширение сидов
- Объединение пересекающих выравниваний
- Построение итогового выравнивания



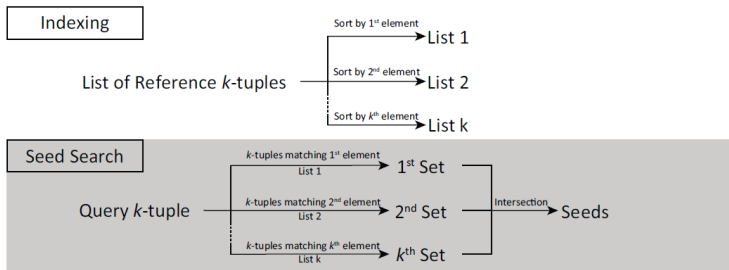
# OMBlast: Поиск стартовых сидов - индексация

Фрагмент  $q$  на карте совпадает с фрагментом  $r$  на референсе:

$$r(1 - T_s) - T_m \leq q \leq r(1 + T_s) + T_m$$

$T_s$  - ошибка масштабирования

$T_m$  - ошибка измерений



# OMBlast: Поиск стартовых сидов - бины

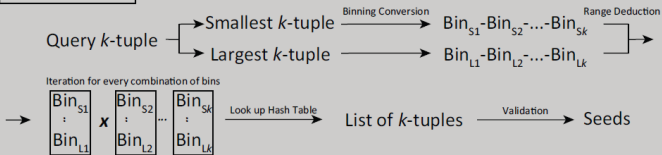
## Binning

Bin	Interval
A	1-5000
B	5001-10000
C	10001-15000
...	...

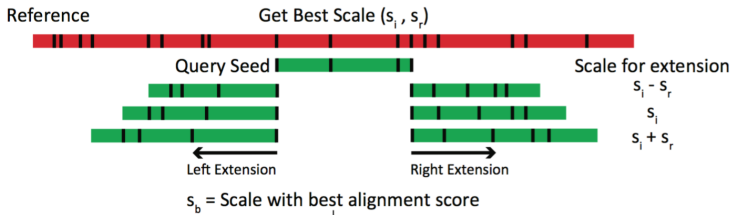
## Indexing

Reference  $k$ -tuple  $\xrightarrow{\text{Binning Conversion}}$   $\text{Bin}_1\text{-Bin}_2\text{-...-Bin}_k \xrightarrow{\text{Hashing}}$  Hash Table

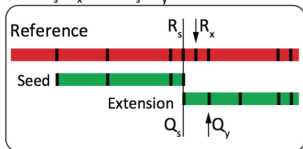
## Seed Search



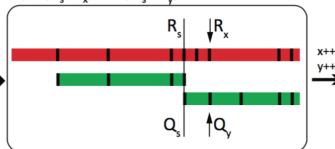
# OMBlast: Расширение сидов



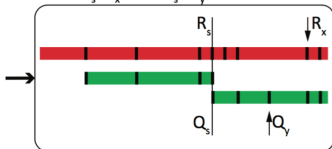
(i)  $D(R_s, R_x) < D(Q_s, Q_y)$



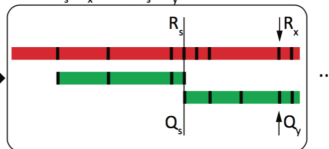
(ii)  $D(R_s, R_x) \sim D(Q_s, Q_y)$



(iii)  $D(R_s, R_x) > D(Q_s, Q_y)$



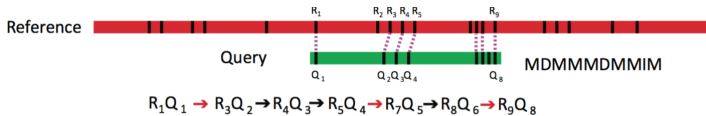
(iv)  $D(R_s, R_x) \sim D(Q_s, Q_y)$



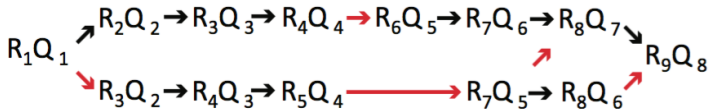
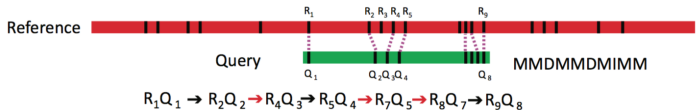
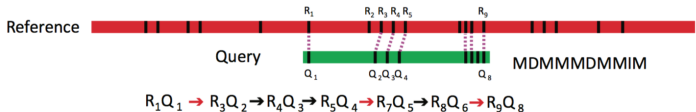
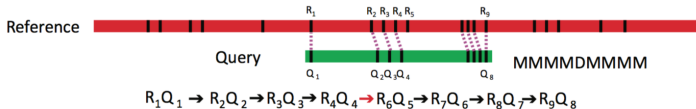
# OMBlast: Объединение выравниваний (1)

Строится взвешенный ациклический граф:

- Вершины - выравненные разрезы
- Рёбра - между двумя парами последовательно (на одной карте) выравненных разрезов
- Веса -  $t_m u_m - t_{es} u_{es} - t_{ms} u_{ms}$   
 $u_m$  - количество совпадений  
 $u_{es}$  - количество лишних разрезов  
 $u_{ms}$  - количество пропущенных разрезов

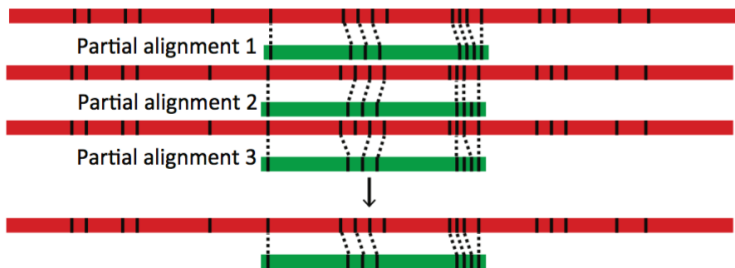


# OMBlast: Объединение выравниваний (2)



# OMBlast: Объединение выравниваний (3)

С помощью динамического программирования определяется путь в графе с наибольшим весом



# OMBlast: Построение итогового выравнивания

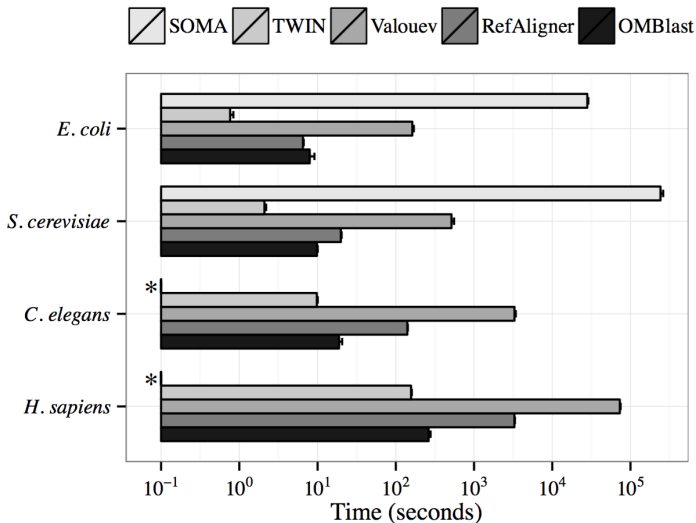
# OMBlast: Результаты - входные данные

Organism	Genome Size (Mbp)	Total Signals	Average Bases Between Signals (kbp)
<i>E. coli</i>	4.6	683	$6.8 \pm 7.3$
<i>S. cerevisiae</i>	12.1	1953	$6.2 \pm 6.7$
<i>C. elegans</i>	100.3	14837	$6.8 \pm 8.0$
<i>H. sapiens</i>	3088.3	377143	$8.2 \pm 83.2$

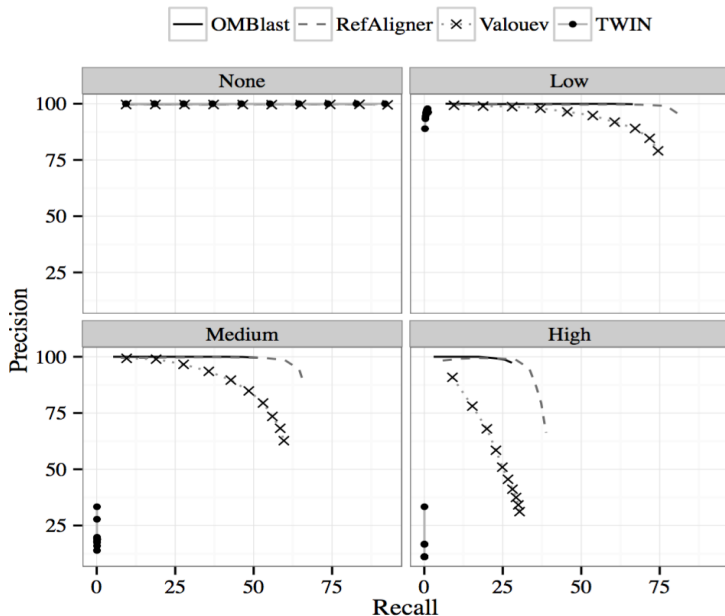
Error Rate	None	Low	Medium	High
Extra Signal Rate	0	0.000005	0.00001	0.00002
Missing Signal Rate	0	0.05	0.1	0.2
Scaling	0	0.02	0.04	0.08
Measurement (bp)	0	500	500	500
Resolution (bp)	0	1200	1200	1200



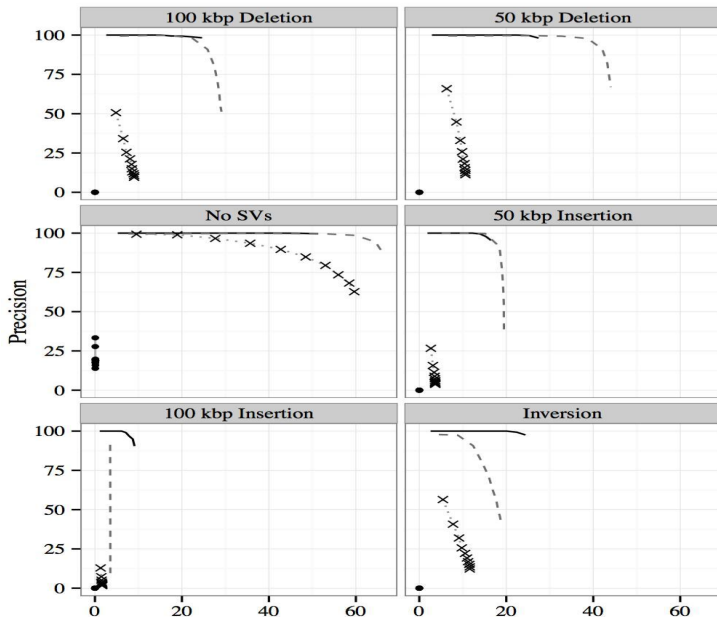
# OMBlast: Результаты - время работы



# OMBlast: Результаты - точность и полнота



# OMBlast: Результаты - наличие SV



Спасибо за внимание!