

# Сборка генома с помощью облаков ридов

Толстогоганов Иван Николаевич

научный руководитель: Банкевич Антон Викторович

СПб АУ НОЦНТ РАН

15 июня 2017 г.

- Поиск новых генов
- Обнаружение сложных вариаций
- Предсказание функции белка

# Облака ридов: общее описание

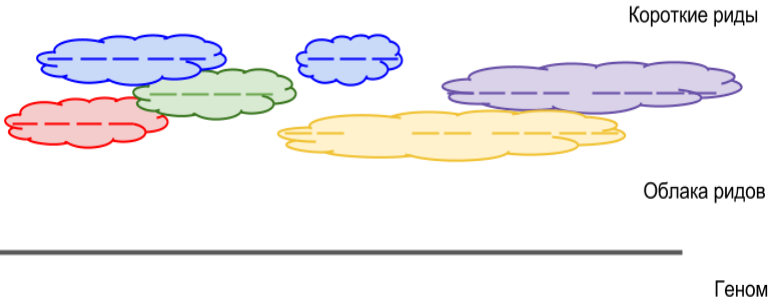
Фрагменты



Геном

- ДНК разделяется на длинные фрагменты (до 50 килобаз)

# Облака ридов: общее описание



- ДНК разделяется на длинные фрагменты (до 50 килобаз)
- Фрагменты секвенируются и баркодируются
- Из одного фрагмента формируется облако коротких ридов

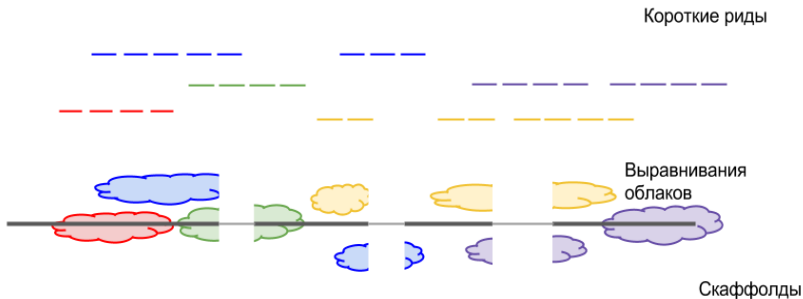
## Цель

Разработка алгоритма для сборки метагенома с помощью облаков ридов

## Задачи

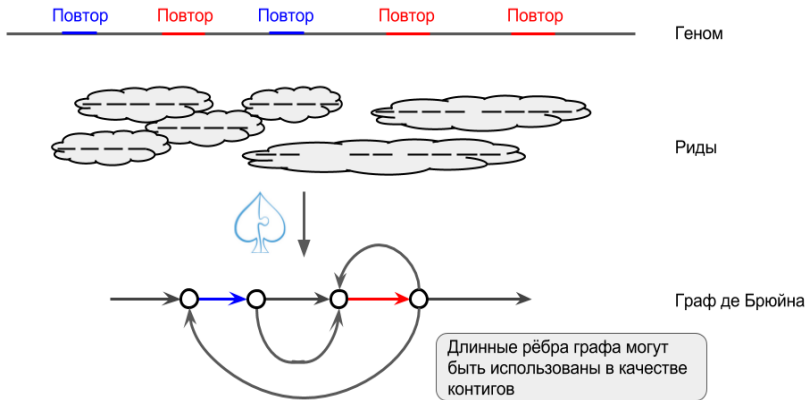
- Реализовать алгоритм в качестве дополнительного модуля к ассемблеру SPAdes
- Поддерживать два наиболее распространённых протокола создания облаков ридов
- Сравнить эффективность алгоритма с альтернативными подходами на экспериментальных данных

# Облака ридов: скаффолдинг

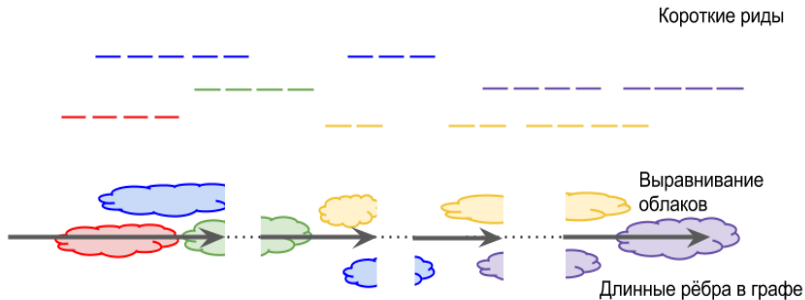


- Облако может быть восстановлено путем выравнивания на контиг
- Выравнивание используется для скаффолдинга

# Граф де Брюйна



# Облака ридов на графе де Брюйна



- Облака могут быть приложены и к рёбрам графа
- Для каждого длинного ребра мы хотим найти следующее в геномном пути
- Для этого используются облака ридов и структура графа





- Длинные рёбра в геномном пути могут быть расположены далеко друг от друга
- Пробелы могут быть закрыты после нахождения правильного порядка



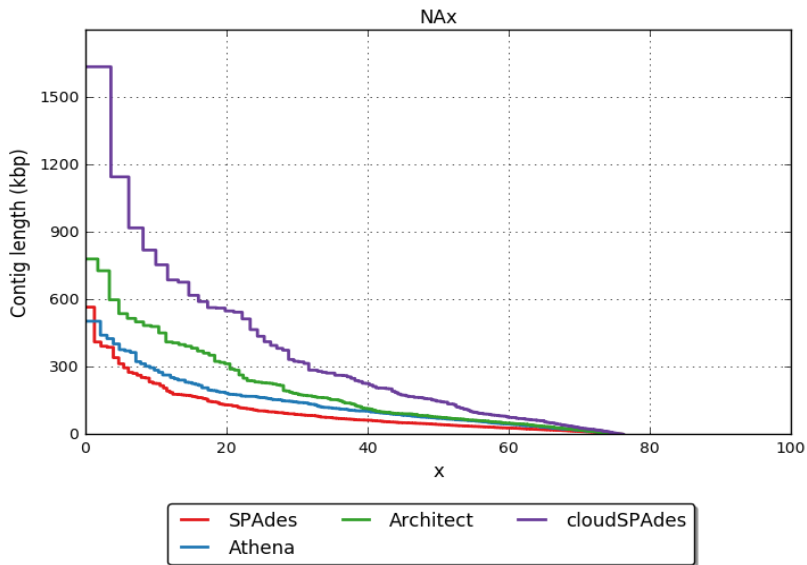
- Длинные рёбра в геномном пути могут быть расположены далеко друг от друга
- Пробелы могут быть закрыты после нахождения правильного порядка
- Для закрытия используются пути из коротких ребер

- Используется простой метагеном из 10 организмов
- Референсы недоступны, вместо них используются:
  - 23 последовательности рибосомальных оперонов
  - Отдельные сборки каждого организма

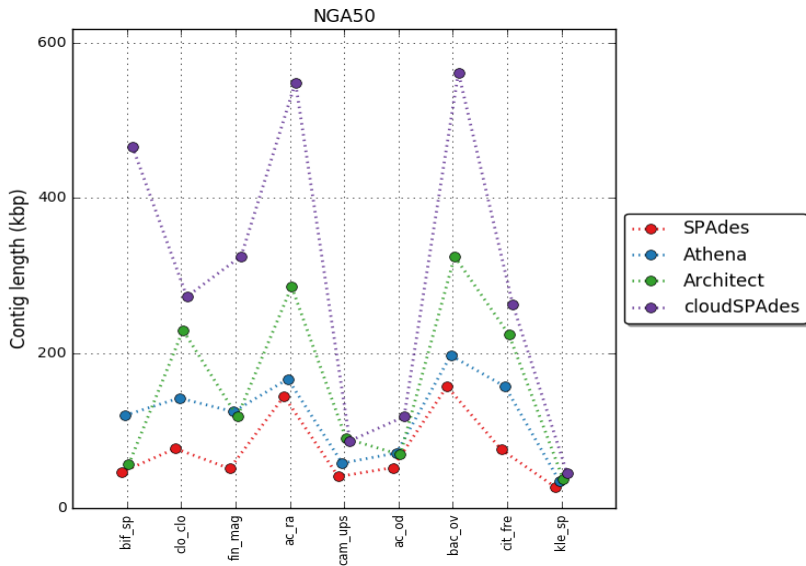
Сборки были переданы в качестве референсов инструменту metaQUAST, предназначенному для оценки качества метагеномной сборки

Assembly	SPAdes	Athena	Architect	cloudSPAdes
# contigs	4333	3700	3989	3607
Total length	44454653	<b>47352283</b>	44467866	45282462
N50	70276	114821	153964	<b>270900</b>
# misassemblies	<b>21</b>	25	29	39
Genome fraction (%)	98.744	98.923	98.778	<b>99.061</b>
# N's per 100 kbp	<b>0.00</b>	<b>0.00</b>	?	127.14
Largest alignment	565835	503302	780028	<b>1634993</b>
NA50	44044	71261	76715	<b>147448</b>

- cloudSPAdes: Реализация алгоритма сборки с помощью облаков ридов на основе SPAdes
- Athena, Architect: Альтернативные инструменты для метагеномной сборки



# NGA50 для каждого организма



## Преимущества cloudSPAdes:

- Длинные контиги
- Восстановлено 22 оперона из 23, нашлось 9  
НОВЫХ

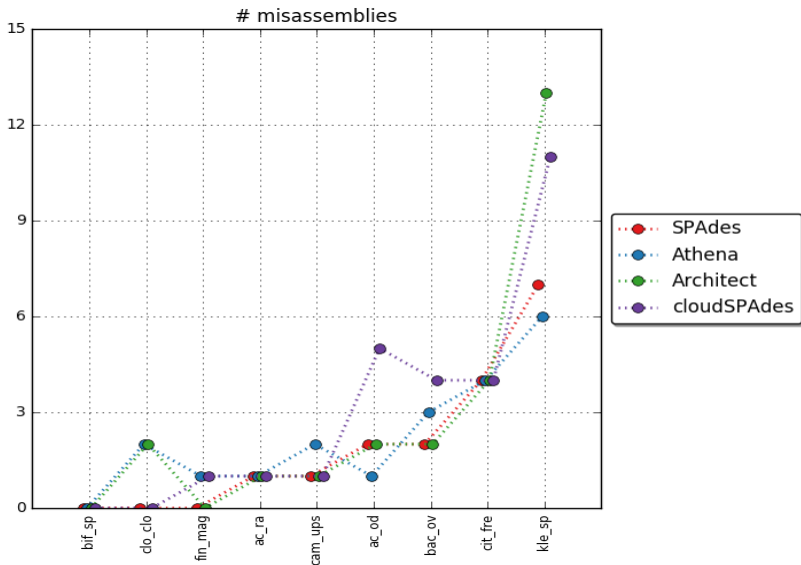
## Недостатки:

- Ошибки сборки в низкопокрытых геномах

- Разработан алгоритм сборки генома с помощью облаков ридов
- Алгоритм встроен в ассемблер SPAdes в качестве дополнительного модуля разрешения повторов
- Поддержаны протоколы создания облаков ридов TruSeq Synthetic Long Reads и 10XGenomics GemCode



# Ошибки сборки



# Полные результаты на GemCode

Assembly	SPAdes	cloudSPAdes	Architect	ARCS	Athena	FragScaff
# contigs	4333	3607	3989	4234	3700	4704
Largest contig	565922	<b>1635422</b>	1163531	565922	506140	1064532
Total length	44454653	45282462	44467866	44455643	47352283	<b>60294921</b>
N50	70276	<b>270900</b>	153964	87970	114821	87521
# misassemblies	<b>21</b>	40	29	62	27	365
Genome fraction (%)	98.745	<b>99.061</b>	98.776	98.750	98.964	98.839
# N's per 100 kbp	<b>0.00</b>	127.14	?	2.23	<b>0.00</b>	?
Largest alignment	565835	<b>1634993</b>	780028	565835	503302	618353
NA50	44044	<b>147327</b>	76715	47089	71261	18146

# Результаты на TSLR

Assembly	Baseline	cloudSPAdes	Architect	FragScaff
# contigs	3657	3372	3407	3293
Total length	46637846	46894699	46683828	<b>52366979</b>
N50	78631	143747	<b>191269</b>	100994
# misassemblies	<b>57</b>	75	186	305
Genome fraction (%)	42.747	42.727	<b>42.786</b>	42.765
# N's per 100 kbp	<b>9.40</b>	624.58	18.31	10755.30
Largest alignment	739469	<b>1311221</b>	1240613	739469
NA50	77174	122883	<b>126089</b>	59381