



Текстовая расшифровка видео:

## СЕКРЕТНЫЕ ОПТИМИЗАЦИИ

**План:**

- Combiner, Partitioner, Comparator, Compression;
- Количество мапперов и редюсеров.

### Combiner, Partitioner, Comparator, Compression

**Combiner** – промежуточный код, который запускается между маппером и редюсером. Может агрегировать данные перед пересылкой их на редюсер.

**Comparator** – возможность указать более точно механизм сравнения ключей при шаффле и сортировке.

Пример того, как запускать job'y с comparator:

```
~$ yarn jar $HADOOP_STREAMING_JAR \  
-D stream.num.map.output.key.fields=2 \  
-D stream.num.reduce.output.key.fields=2 \  
-D mapreduce.job.output.key.comparator.class=org.apache.hadoop.mapreduce.lib.partition.  
-D mapreduce.partition.keycomparator.options="-k1,1n -k2,2" \  
^C
```

Наши данные могут быть нетривиальной структуры, и мы можем захотеть их отсортировать, например, по нестандартному ключу и т.д. Поэтому, когда мы используем `hadoop streaming`, для этих целей напрямую написать `comparator` на Python мы не сможем, но мы сможем указать класс Java и аргументы, которые будем использовать в `comparator`.

В примере в параметрах указан ключ и у маппера, и после маппера, и на входе редюсера может иметь два поля. Далее, для сравнения значений между собой использован класс Java – `KeyFieldBasedComparator` как для ключей, составных из двух полей.



