



# **Сжатие данных в HDFS и YARN, горячие и холодные данные**

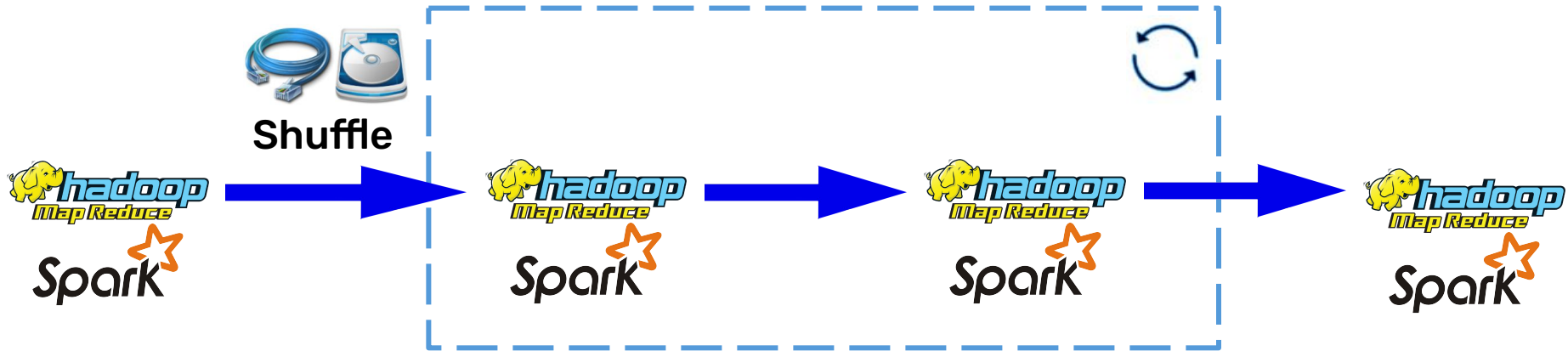
**Драль Алексей**, [study@bigdatateam.org](mailto:study@bigdatateam.org)

CEO at BigData Team, <https://bigdatateam.org>

<https://www.facebook.com/bigdatateam>



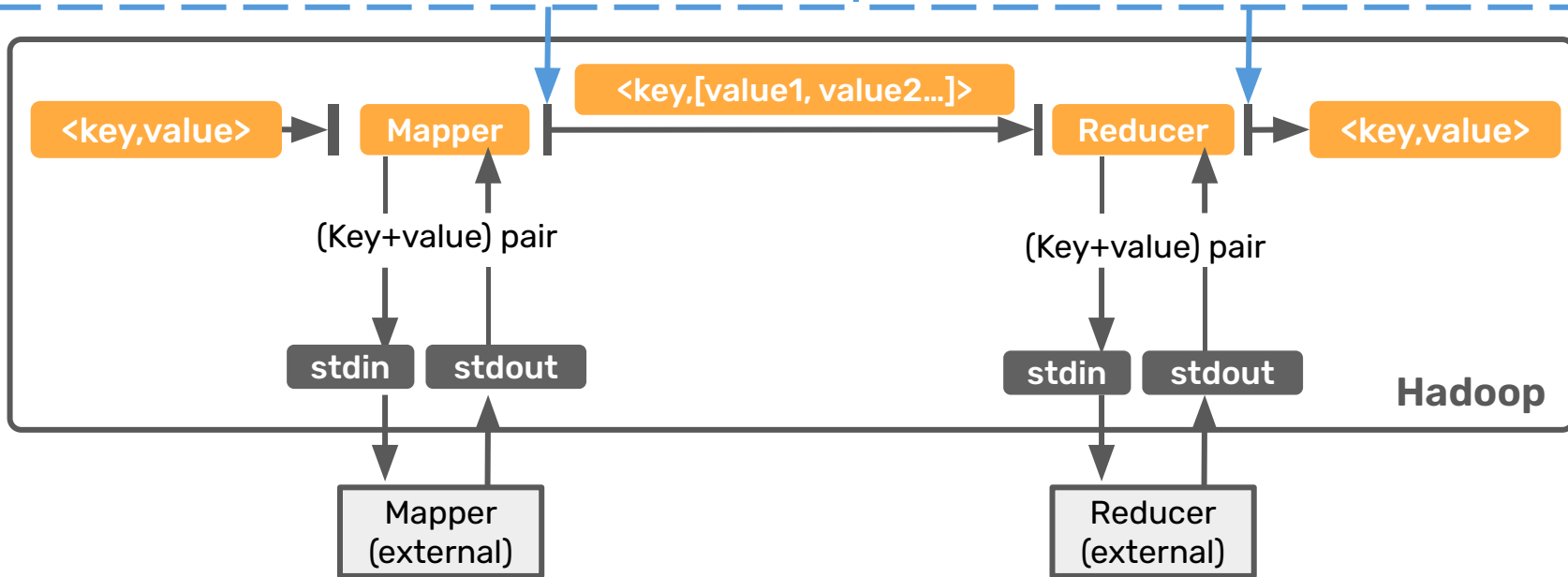
# Disk I/O, Network Bandwidth





# Сжатие данных в Hadoop

```
-D mapreduce.compress.map.output=true  
-D mapreduce.map.output=compression.codec=...  
-D mapreduce.output.compress=true  
-D mapreduce.output=compression.codec=...
```



Spark: `spark.io.compression.codec=...`, `spark.shuffle.compress=true`



**BIGDATA  
TEAM**

# Ликбез по кодекам



- ▶ Кодировка (encoding)



- ▶ Кодировка (encoding)
- ▶ Шифрование (encryption)



- ▶ Кодировка (encoding)
- ▶ Шифрование (encryption)
- ▶ Сжатие (compression)



- ▶ Кодировка (encoding)
- ▶ Шифрование (encryption)
- ▶ Сжатие (compression)

Кодек = encoding + encryption + compression





- ▶ Кодировка (encoding)
- ▶ Шифрование (encryption)
- ▶ Сжатие (compression)

Кодек = encoding + encryption + compression

Hadoop codec = compression



# Compression Codecs

Формат сжатия	Hadoop CompressionCodec
DEFLATE	org.apache.hadoop.io.compress.DefaultCodec
gzip	org.apache.hadoop.io.compress.GzipCodec
bzip2	org.apache.hadoop.io.compress.BZip2Codec
LZO	com.hadoop.compression.lzo.LzopCodec
LZ4	org.apache.hadoop.io.compress.Lz4Codec
Snappy	org.apache.hadoop.io.compress.SnappyCodec



**BIGDATA  
TEAM**

# Trade-off сжатия данных





# Стандартные алгоритмы сжатия

Формат сжатия	Splittable	Комментарии
.deflate .gz (gzip)	Нет	Используется алгоритм DEFLATE



# Стандартные алгоритмы сжатия

Формат сжатия	Splittable	Комментарии
.deflate .gz (gzip)	Нет	Используется алгоритм DEFLATE
.bz2 (bzip)	Да	Более эффективен чем gzip, но сжатие медленнее



# Стандартные алгоритмы сжатия

Формат сжатия	Splittable	Комментарии
.deflate .gz (gzip)	Нет	Используется алгоритм DEFLATE
.bz2 (bzip)	Да	Более эффективен чем gzip, но сжатие медленнее
.lzo	Да*	Распаковка быстрее чем gzip, сжатие менее эффективно

\*можно сделать splittable с помощью добавления индексов



# Стандартные алгоритмы сжатия

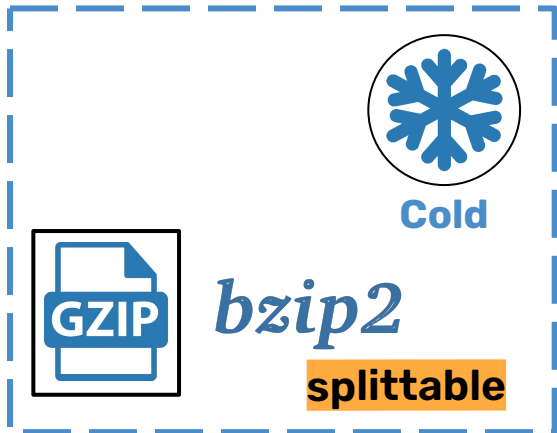
Формат сжатия	Splittable	Комментарии
.deflate .gz (gzip)	Нет	используется алгоритм DEFLATE
.bz2 (bzip)	Да	более эффективен чем gzip, но сжатие медленнее
.lzo	Да*	распаковка быстрее чем gzip, сжатие менее эффективно
.lz4 .snappy	Нет	{у,рас}паковка быстрее чем LZ0

\*можно сделать splittable с помощью добавления индексов



**BIGDATA**  
**TEAM**

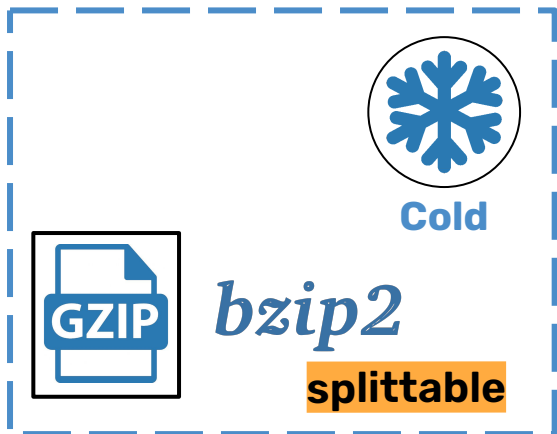
# Стандартные подходы







# Стандартные подходы





# Стандартные подходы

