

2026

Yaroslav Tkachenko

# MAKING KAFKA SERIALIZATION FAST AGAIN: STOP TRUSTING DEFAULTS



Yaroslav Tkachenko

## My background

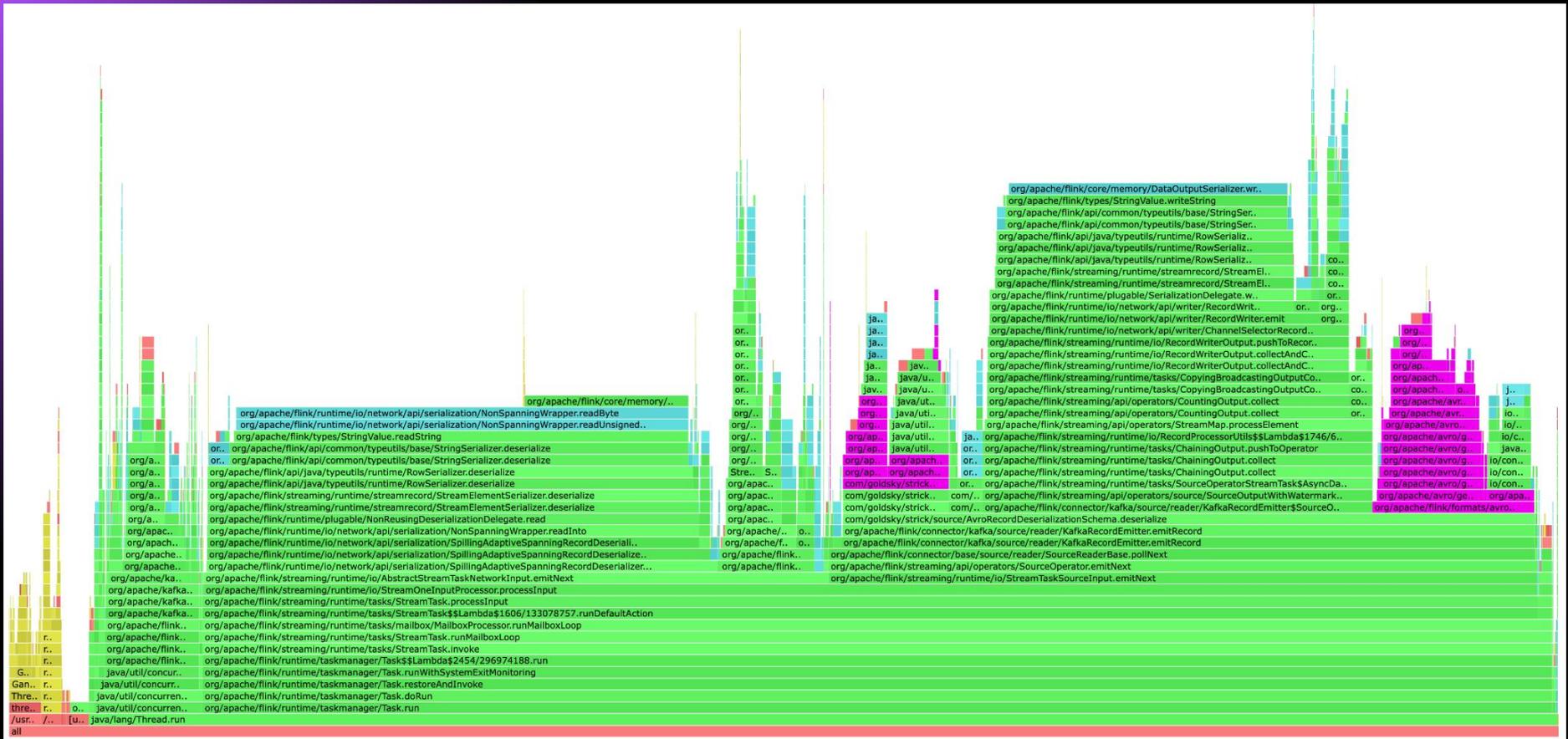
- Started my career as a full-stack web developer, worked in a few startups
- Technical Lead at **Activision**, **Shopify**
- Founding Engineer at **Goldsky**
- Data streaming tooling, consulting and training at **Irontools**

<https://www.linkedin.com/in/sap1ens/>

Your data serialization is slow  
And likely can be 2x - 3x faster

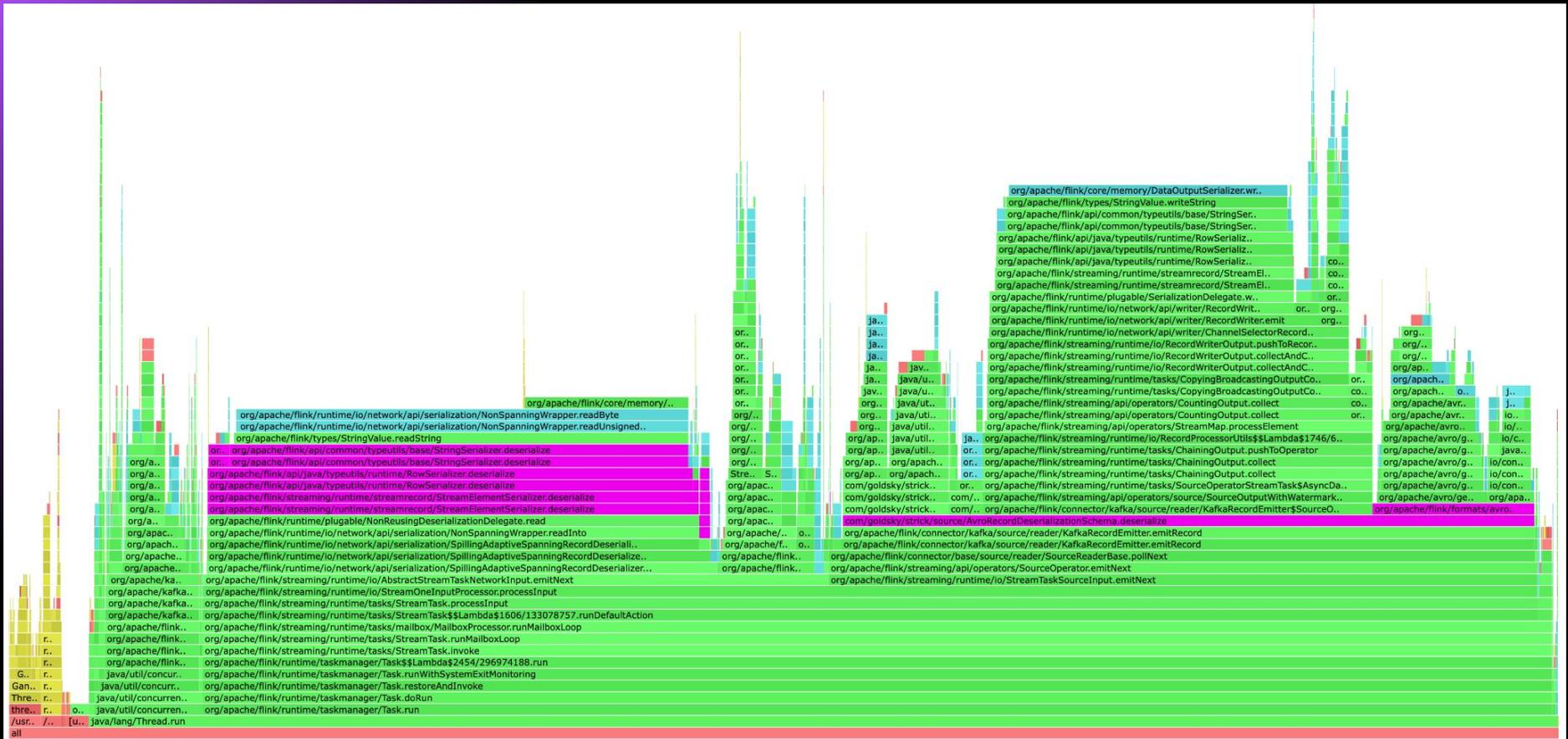
OPINION

Serialization speed can be the biggest contributing factor to the performance of **stateless** streaming pipelines



Matched: 17.28% X

## Kafka Avro Flink Source Deserializer



```

org/apache/flink/core/memory/...
org/apache/flink/runtime/io/network/api/serialization/NonSpanningWrapper.readByte
org/apache/flink/runtime/io/network/api/serialization/NonSpanningWrapper.readUnsigned...
org/apache/flink/types/StringValue.readString
org/apache/flink/api/common/typeutils/base/StringSerializer.deserialize
org/apache/flink/api/common/typeutils/base/StringSerializer.deserialize
org/apache/flink/api/java/typeutils/runtime/RowSerializer.deserialize
org/apache/flink/api/java/typeutils/runtime/RowSerializer.deserialize
org/apache/flink/streaming/runtime/streamrecord/StreamElementSerializer.deserialize
org/apache/flink/streaming/runtime/streamrecord/StreamElementSerializer.deserialize
org/apache/flink/runtime/plugable/NonReusingDeserializationDelegate.read
org/apache/flink/runtime/io/network/api/serialization/NonSpanningWrapper.readInto
org/apache/flink/runtime/io/network/api/serialization/SpillingAdaptiveSpanningRecordDeseriali...
org/apache/flink/runtime/io/network/api/serialization/SpillingAdaptiveSpanningRecordDeseriali...
org/apache/flink/streaming/runtime/io/AbstractStreamTaskNetworkInput.emitNext
org/apache/flink/streaming/runtime/io/StreamOneInputProcessor.processInput
org/apache/flink/streaming/runtime/tasks/StreamTask.processInput
org/apache/flink/streaming/runtime/tasks/StreamTask$Lambda$1606/133078757.runDefaultAction
org/apache/flink/streaming/runtime/tasks/mailbox/MailboxProcessor.runMailboxLoop
org/apache/flink/streaming/runtime/tasks/StreamTask.runMailboxLoop
org/apache/flink/streaming/runtime/tasks/StreamTask.invoke
org/apache/flink/runtime/taskmanager/Task$Lambda$2454/296974188.run
java/util/concurrent/Task.runWithSystemExitMonitoring
org/apache/flink/runtime/taskmanager/Task.restoreAndInvoke
org/apache/flink/runtime/taskmanager/Task.doRun
java/util/concurrent/Task.run
java/lang/Thread.run

```

Matched: 77.03% X

## Kafka Avro Flink Source Deserializer





# THINGS CAN GET WORSE

→ Complex schemas

- Arrays, nested objects, etc.

→ Large payloads

- E.g. WarpStream recommends setting `max.request.size` = **64000000**

→ Number of Kafka topic partitions can be a bottleneck

- Unless using non-standard multi-threading approaches for scaling

What can we do?

# KEY OPTIMIZATIONS

## Vectorization / SIMD

Leverage modern CPU instructions to process multiple data points at once.

- Most modern CPUs support it.
- Historically supported well in C++.
- Java is lagging\*

## Code Specialization

Generate specialized code in runtime.

- Highly specialized code can be very compact, eliminate branching and function calls.
- JIT-friendly.
- Works great with schemas and schema registries.

## Columnar Formats

Organizing data in columns.

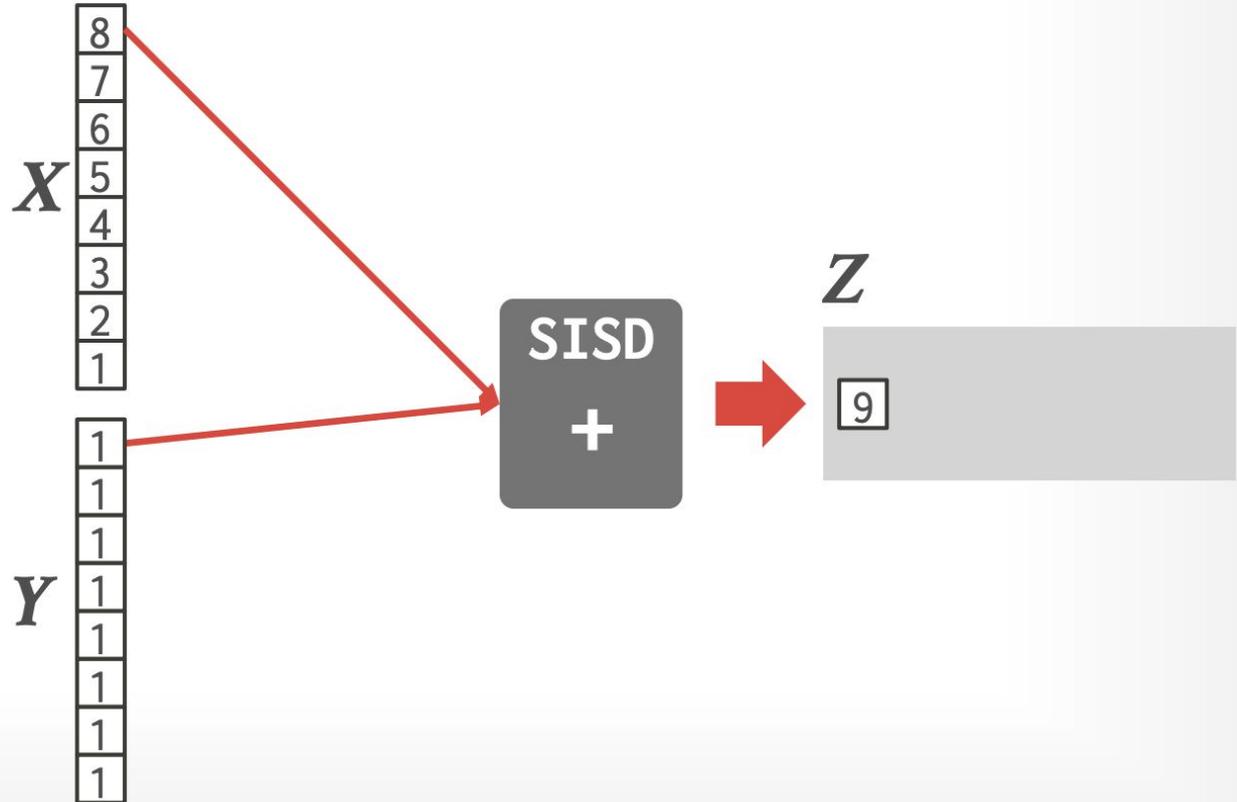
- Columnar storage and processing have been used for many years in OLAP systems.
- Enables CPU cache locality.
- Enables efficient vectorization.

# Vectorization / SIMD

$$X + Y = Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

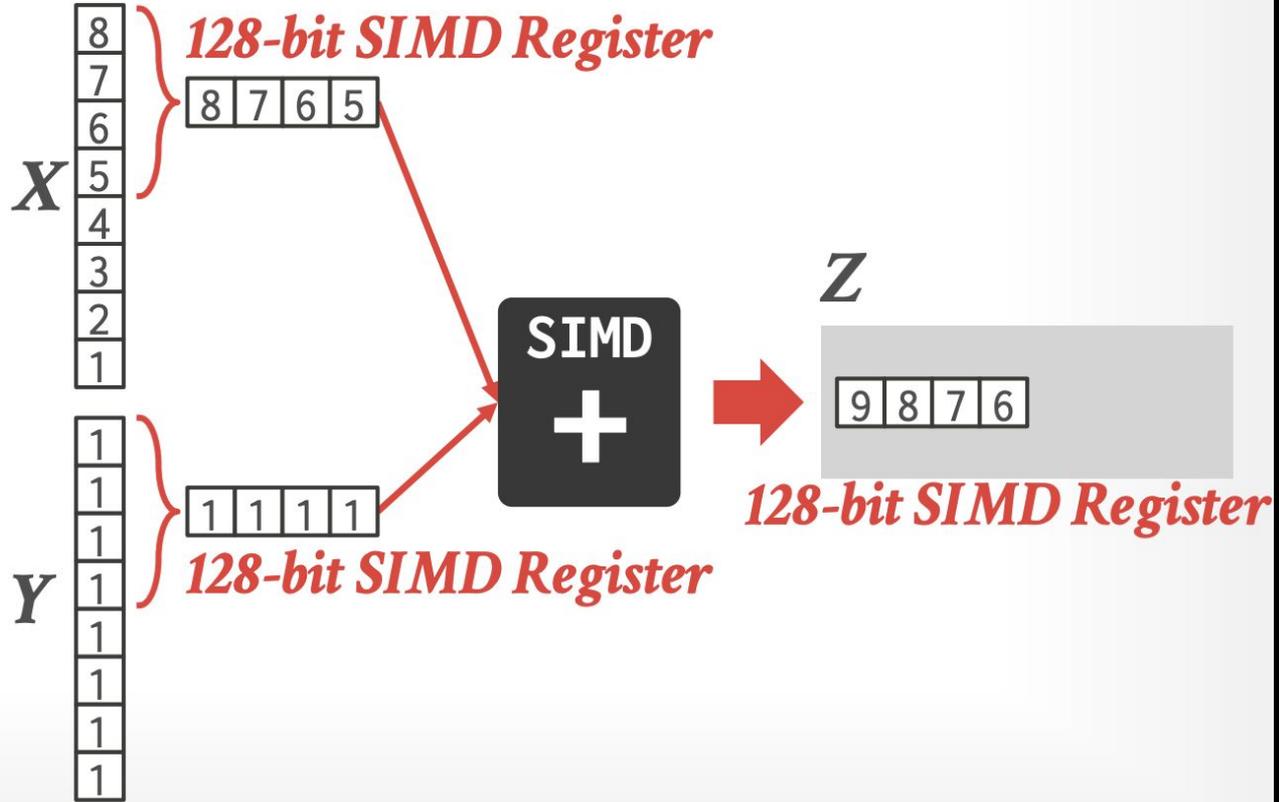
```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



$$X + Y = Z$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

```
for (i=0; i<n; i++) {  
    Z[i] = X[i] + Y[i];  
}
```



# MANUAL VECTORIZATION

Use CPU intrinsics to manually marshal data between SIMD registers and execute vectorized instructions.

Probably not something most people enjoy. Look for libraries that implement this support.

```
void add(int *X,  
         int *Y,  
         int *Z) {  
    __mm128i *vecX = (__m128i*)X;  
    __mm128i *vecY = (__m128i*)Y;  
    __mm128i *vecZ = (__m128i*)Z;  
    for (int i=0; i<MAX/4; i++) {  
        _mm_store_si128(vecZ++,  
            ↪ _mm_add_epi32(*vecX++,  
                ↪ *vecY++));  
    }  
}
```

Java Version	Status	Package	ByteVector availability
15	<i>Preview only</i> (very early)	<code>jdk.incubator.vector</code>	✗ Not usable yet (incomplete)
16	<b>1st incubator (JEP 338)</b>	<code>jdk.incubator.vector</code>	✓ Available
17 (LTS)	2nd incubator (JEP 414)	<code>jdk.incubator.vector</code>	✓ Available
18	3rd incubator (JEP 417)	<code>jdk.incubator.vector</code>	✓ Available
19	4th incubator (JEP 426)	<code>jdk.incubator.vector</code>	✓ Available
20	5th incubator (JEP 438)	<code>jdk.incubator.vector</code>	✓ Available
21 (LTS)	<b>Finalized Vector API</b>	<code>jdk.vector</code> (not incubator)	⚠ Package moved
22+	Final	<code>jdk.vector</code>	

# SIMD vs SWAR

## SIMD

```
var vec =  
ByteVector.fromArray(SPECIES_256,  
array, i);  
  
var cmp = vec.eq((byte)0);  
  
if (cmp.anyTrue()) {  
    // zero found  
}
```

## SWAR

```
long x = load64(ptr);  
  
long tmp = (x -  
0x0101010101010101L) & ~x &  
0x8080808080808080L;  
  
if (tmp != 0) {  
    // zero byte present  
}
```

# Code Specialization

Can we perform the **same** amount of work with **fewer** CPU instructions?

```
Object read(Object obj, Schema schema) {  
    switch(schema.type()) {  
        case STRUCT:  
            return readStruct(obj, schema);  
        case STRING:  
            return readString(obj, schema);  
        // ... rest of the types  
    }  
}
```

```
Object readStruct(Object obj, Schema schema) {  
    Struct struct = new Struct();  
    for (Field field: schema.fields()) {  
        struct.field(field.name(), read(field.value(), field.schema()));  
    }  
    return struct;  
}
```

“Interpreted” Approach

```
User read(Object obj, Schema schema) {
    User user = new User();
    Decoder decoder = new Decoder(obj, schema);

    user.put(0, decoder.readString());
    user.put(1, decoder.readString());
    user.put(2, decoder.readLong());

    Order order = new Order();
    order.put(0, decoder.readString());
    user.put(3, order);

    return user;
}
```

# RUNTIME CODE SPECIALIZATION IN JVM

→ You can generate Java code, compile and instantiate it with a classloader. So you'll need:

- Code generation: via some DSL or templating
- Compiler: Janino, javac (!)
- A way to load / switch implementations

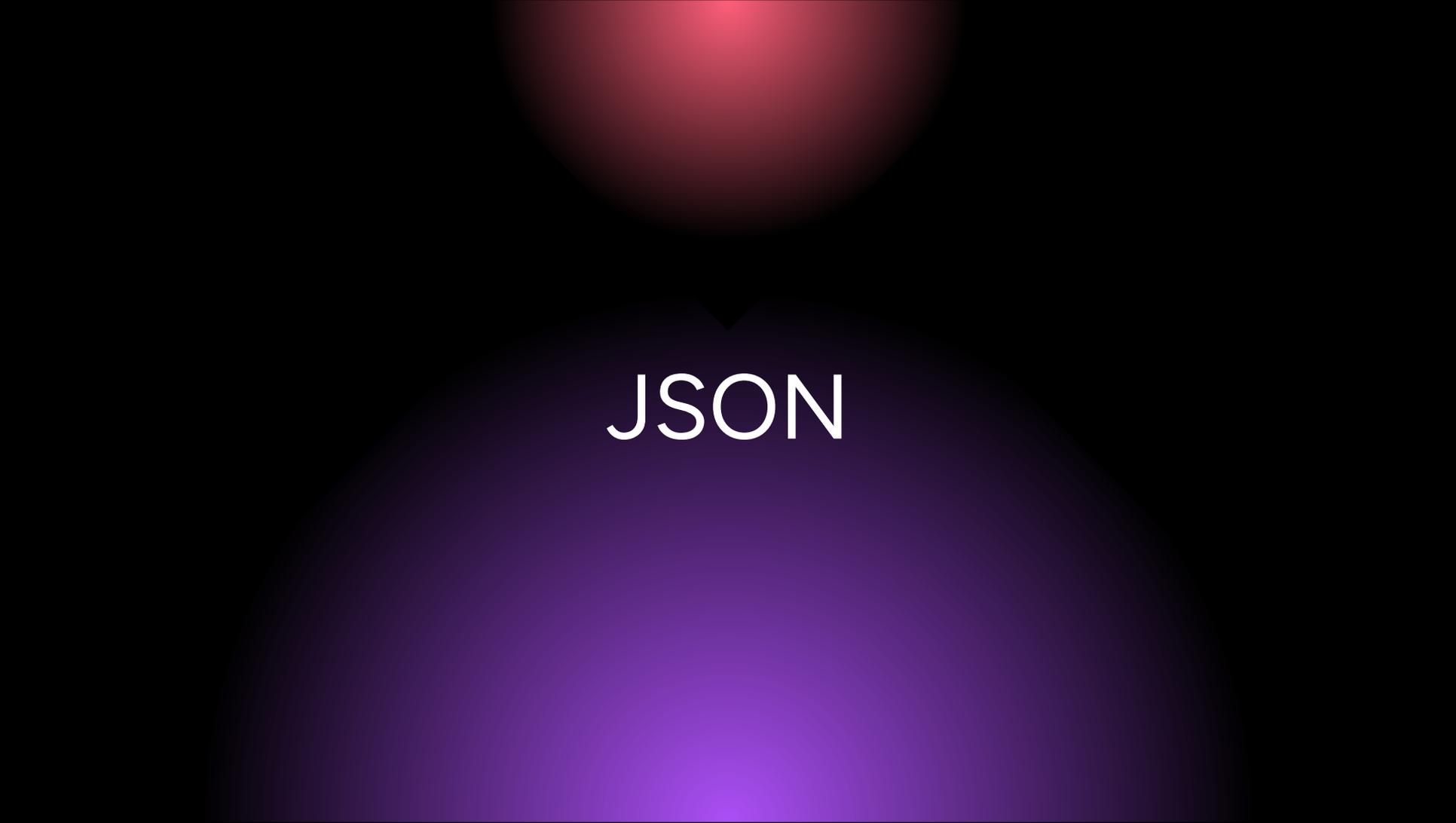
→ Or you can generate bytecode directly

- ASM is a very popular library for that

# CODE SPECIALIZATION FOR DESERIALIZATION

For every new detected schema:

- Create a specialized deserializer tailored to the schema
- Use that instead of the default serializer
  - Can still fallback to the default if needed
  - If compilation takes a while, use the default until it's done, then switch



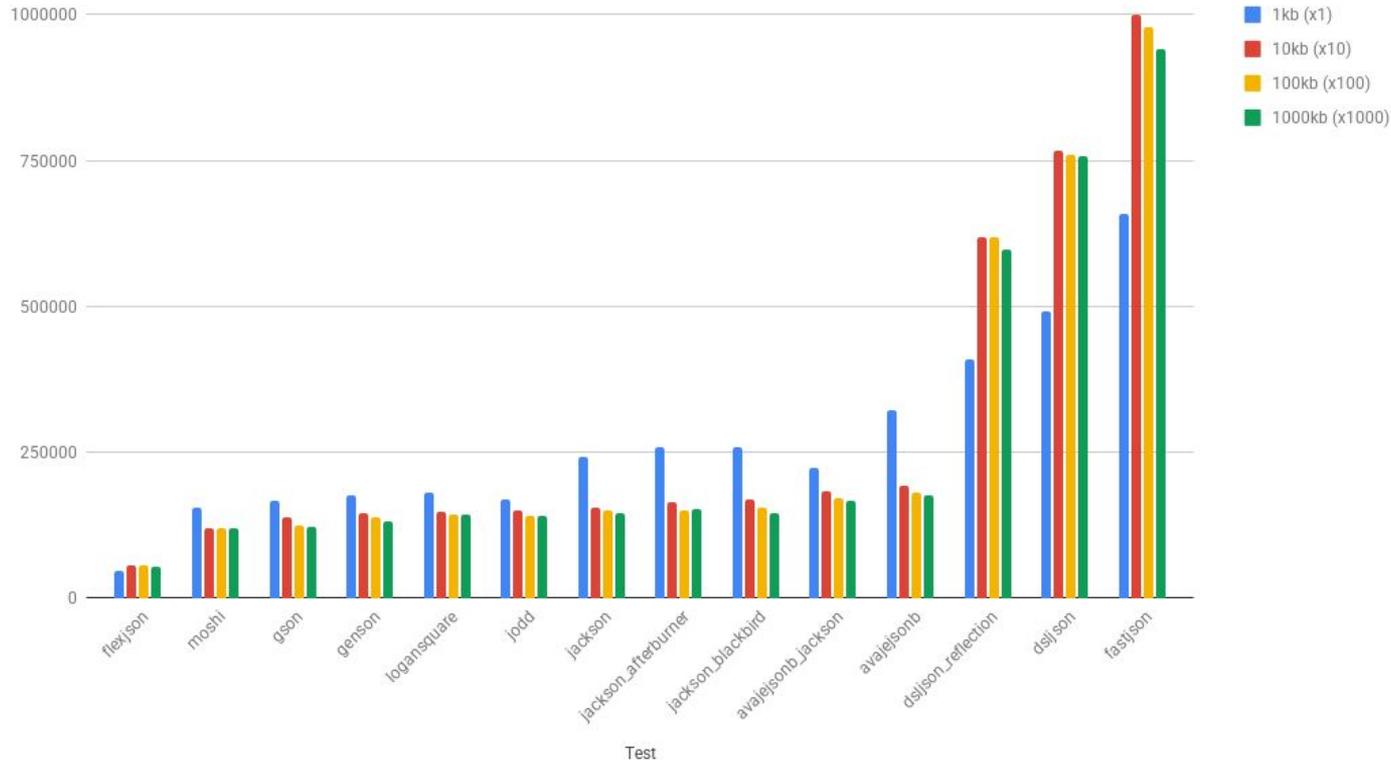
JSON

# JSON DESERIALIZATION

When it comes to Java JSON serialization / deserialization, **Jackson** is everywhere

- Confluent's KafkaJsonDeserializer
- Spring's JsonDeserializer
- Flink's JsonDeserializationSchema

## Deserialization of various payload sizes - ops/s (higher is better)



# simdjson-java

CI passing License Apache 2

A Java version of [simdjson](#) - a JSON parser using SIMD instructions, based on the paper [Parsing Gigabytes of JSON per Second](#) by Geoff Langdale and Daniel Lemire.

## Code Sample

### DOM Parser

```
byte[] json = loadTwitterJson();

SimdJsonParser parser = new SimdJsonParser();
JsonValue jsonValue = parser.parse(json, json.length);
Iterator<JsonValue> tweets = jsonValue.get("statuses").arrayIterator();
while (tweets.hasNext()) {
    JsonValue tweet = tweets.next();
    JsonValue user = tweet.get("user");
    if (user.get("default_profile").asBoolean()) {
        System.out.println(user.get("screen_name").asString());
    }
}
```

### Schema-Based Parser

```
byte[] json = loadTwitterJson();

SimdJsonParser parser = new SimdJsonParser();
SimdJsonTwitter twitter = simdJsonParser.parse(buffer, buffer.length, SimdJsonTwitter.class);
for (SimdJsonStatus status : twitter.statuses()) {
    SimdJsonUser user = status.user();
    if (user.default_profile()) {
        System.out.println(user.screen_name());
    }
}
```

```
1. public class SimdJsonDeserializer<T> implements Deserializer<T> {
2.
3.     private final SimdJsonParser parser = new SimdJsonParser();
4.     private final Class<T> clazz;
5.
6.     public SimdJsonDeserializer(Class<T> clazz) {
7.         this.clazz = clazz;
8.     }
9.
10.    @Override
11.    public T deserialize(String topic, byte[] data) {
12.        if (data == null || data.length == 0) {
13.            return null;
14.        }
15.        return parser.parse(data, data.length, clazz);
16.    }
17. }
```

“We require Java 24 or better”

“IllegalArgumentException:  
Unsupported vector species:  
Species[byte, 16, S\_128\_BIT]”

## BENCHMARK

Benchmark	Mode	Cnt	Score	Error	Units
Benchmark.jackson	avgt	10	565.766 ±	2.474	ns/op
Benchmark.simdJava	avgt	10	8878.502 ±	29.739	ns/op

build passing coverage 74% maven-central v2.0.61.android8 release v2.0.61 Java 8+ license Apache-2.0  
gitpod ready-to-code latest snapshot v2.0.58-SNAPSHOT Stars 4.3k Forks 551 user repos 0 contributors 111

English Documentation | 中文文档

The issues of fastjson will be also posted on [Alibaba Cloud Developer Community](#)

## FASTJSON v2

FASTJSON v2 is an upgrade of the FASTJSON , with the goal of providing a highly optimized JSON library for the next ten years.

- Supports the JSON and JSONB Protocols.
- Supports full parsing and partial parsing.
- Supports Java servers and Android Clients, and has big data applications.
- Supports Kotlin [https://alibaba.github.io/fastjson2/Kotlin/kotlin\\_en](https://alibaba.github.io/fastjson2/Kotlin/kotlin_en)
- Supports Android 8+
- Supports JSON Schema [https://alibaba.github.io/fastjson2/JSONSchema/json\\_schema\\_en](https://alibaba.github.io/fastjson2/JSONSchema/json_schema_en)



fastjson2 uses code specialization and SWAR (and other optimizations)

```
1. public class IronJsonDeserializer<T> implements Deserializer<T> {
2.
3.     private final Class<T> clazz;
4.     private final ObjectReader<T> objectReader;
5.     private final JSONReader.Context jsonContext;
6.
7.     public IronJsonDeserializer(Class<T> clazz) {
8.         this.clazz = clazz;
9.         ObjectReaderProvider provider = JSONFactory.getDefaultObjectReaderProvider();
10.        this.jsonContext = new JSONReader.Context(provider);
11.        this.objectReader = provider.getObjectReader(clazz,false);
12.    }
13.
14.    @Override
15.    public T deserialize(String topic, byte[] data) {
16.        if (data == null || data.length == 0) {
17.            return null;
18.        }
19.
20.        try (JSONReader reader = JSONReader.of(data, jsonContext)) {
21.            return objectReader.readObject(reader, clazz,null, 0);
22.        }
23.    }
24. }
```

```
1. // Create class writer
2. ClassWriter cw = new ClassWriter(ClassWriter.COMPUTE_MAXS | ClassWriter.COMPUTE_FRAMES);
3.
4. // Define class
5. cw.visit(
6.     Opcodes.V1_8,
7.     Opcodes.ACC_PUBLIC,
8.     internalName,
9.     null,
10.    "java/lang/Object",
11.    null
12. );
13.
14. // Add public fields
15. for (int i = 0; i < rowType.getFieldCount(); i++) {
16.     String fieldName = rowType.getFieldNames().get(i);
17.     LogicalType fieldType = rowType.getTypeAt(i);
18.     String typeDescriptor = getTypeDescriptor(fieldType);
19.     String signature = getTypeSignature(fieldType);
20.
21.     // Add public field
22.     FieldVisitor fv = cw.visitField(
23.         Opcodes.ACC_PUBLIC,
24.         fieldName,
25.         typeDescriptor,
26.         signature,
27.         null
28.     );
29.
30.     // Add JSONField annotation to field
31.     AnnotationVisitor av = fv.visitAnnotation("Lcom/alibaba/fastjson2/annotation/JSONField;", true);
32.     av.visit("name", fieldName);
33.     av.visitEnd();
34.     fv.visitEnd();
35. }
```

```
1.  switch (type.getTypeRoot()) {
2.      case CHAR:
3.      case VARCHAR:
4.          return "Ljava/lang/String;";
5.      case BOOLEAN:
6.          return "Ljava/lang/Boolean;";
7.      case TINYINT:
8.          return "Ljava/lang/Byte;";
9.      case SMALLINT:
10.         return "Ljava/lang/Short;";
11.     case INTEGER:
12.         return "Ljava/lang/Integer;";
13.     // ...
14. }
```

# BENCHMARK

Benchmark	Mode	Cnt	Score	Error	Units
Benchmark.fastjson2	avgt	10	173.572 ±	0.913	ns/op
Benchmark.jackson	avgt	10	565.766 ±	2.474	ns/op
Benchmark.simdJava	avgt	10	8878.502 ±	29.739	ns/op



Avro

# AVRO DESERIALIZATION

When it comes to Java Avro serialization / deserialization, almost always the standard **org.apache.avro** library is used.

- Confluent's KafkaAvroDeserializer
- Flink's AvroDeserializationSchema
- ...

# Avro-Util

push flow passing latest 0.4.38

A collection of utilities and libraries to allow java projects to better work with avro.

## How To Use

Artifacts are published to [artifactory](#). To use them (say in a gradle build) you'd need to add that repository to your build.gradle, and then add a dependency on the module(s) you wish to use:

```
repositories {  
  maven {  
    url "https://linkedin.jfrog.io/artifactory/avro-util/"  
  }  
}  
...  
dependencies {  
  compile "com.linkedin.avroutil:helper-all:<latest version>"  
}
```



## Background

Apache Avro is a widely used serialization format.

Unfortunately it is not always possible to write java code that "just works" across multiple versions of avro, as there have been breaking changes to both the API and wire format. This project provides utility code to enable java developers to write code that is compatible across a wide range of avro versions

Furthermore, this project enables Avro-based projects to achieve better performance.

avro-fastserde uses code specialization, it starts compilation for a new schema asynchronously and performs a swap when it's ready

```
[INFO ] 2026-02-28 15:45:56.586 [avro-fastserde-compile-thread-1] FastSerdeBase - Starting compilation for the generated source file: User_GenericDeserializer_739236478_739236478.java
```

```
[INFO ] 2026-02-28 15:45:56.868 [avro-fastserde-compile-thread-1] FastSerdeBase - Successfully compiled class User_GenericDeserializer_739236478_739236478 defined at source file User_GenericDeserializer_739236478_739236478.java
```

```
[DEBUG] 2026-02-28 15:45:56.870 [avro-fastserde-compile-thread-1] FastSerdeCache - Generated classes dir: /var/folders/jl/60lw5h01623gd51ffm22hzrr0000gn/T/generated8742529887734165862 and generation of generic FastDeserializer is done for writer schema of type: example.avro.User with fingerprint: {  
  "type" : "record",  
  "name" : "User",  
  "namespace" : "example.avro",  
  "fields" : [ {  
    "name" : "id",  
    "type" : "int"  
  }  
},
```

```
[DEBUG] 2026-02-28 15:45:57.553 [main] FastGenericDatumReader - FastGenericDeserializer was generated and cached for reader schema:
```

```
[{"type":"record","name":"User","namespace":"example.avro","fields":[{"name":"id","type":"int"},
```

```
1. public IndexedRecord deserialize(IndexedRecord reuse, Decoder decoder,
DatumReaderCustomization customization) throws IOException {
2.     return deserializeUser0(null, (decoder), (customization));
3. }
4.
5. public IndexedRecord deserializeUser0(Object reuse, Decoder decoder,
DatumReaderCustomization customization) throws IOException {
6.     IndexedRecord User;
7.     User = new org.apache.avro.generic.GenericData.Record(readerSchema);
8.     User.put(0, (decoder.readInt()));
9.     Utf8 charSequence0 = (decoder).readString(null);
10.    User.put(1, charSequence0);
11.    Utf8 charSequence1 = (decoder).readString(null);
12.    User.put(2, charSequence1);
13.    // ..
14. }
```

# BENCHMARK

Benchmark	Mode	Cnt	Score	Error	Units
Benchmark.avroDeserializationComplex	avgt	10	610.802 ±	10.620	ns/op
Benchmark.avroDeserializationStandard	avgt	10	384.341 ±	4.519	ns/op
Benchmark.avroDeserializationTrivial	avgt	10	131.794 ±	4.027	ns/op
Benchmark.ironAvroDeserializationComplex	avgt	10	224.669 ±	7.045	ns/op
Benchmark.ironAvroDeserializationStandard	avgt	10	134.520 ±	3.171	ns/op
Benchmark.ironAvroDeserializationTrivial	avgt	10	44.185 ±	0.609	ns/op

# SUMMARY

## VECTORIZATION / SIMD

- Leverage modern CPU instructions to process multiple data points at once.
- Check various SIMD JSON libraries, e.g. simdjson-java.
- fastjson2 uses SIMD within a register (SWAR) technique, not specialized SIMD instructions\*.

## CODE SPECIALIZATION / COMPILATION

- Specialized compiled code with minimum instructions is always faster.
- Both fastjson2 and avro-fastserde compile specialized serializer / deserializer classes for a given schema (or class) in runtime.
- Still compatible with Schema Registries.

OPINION

Most importantly, do not trust the defaults: measure things yourself.

# QUESTIONS?

NEWSLETTER!

<https://streamingdata.tech>

