



CODE GENERATION WITH PROTOBUF @ CANVA

Wade Jensen

Arthur Baudry



@wadejensen



@ArthurBaudry

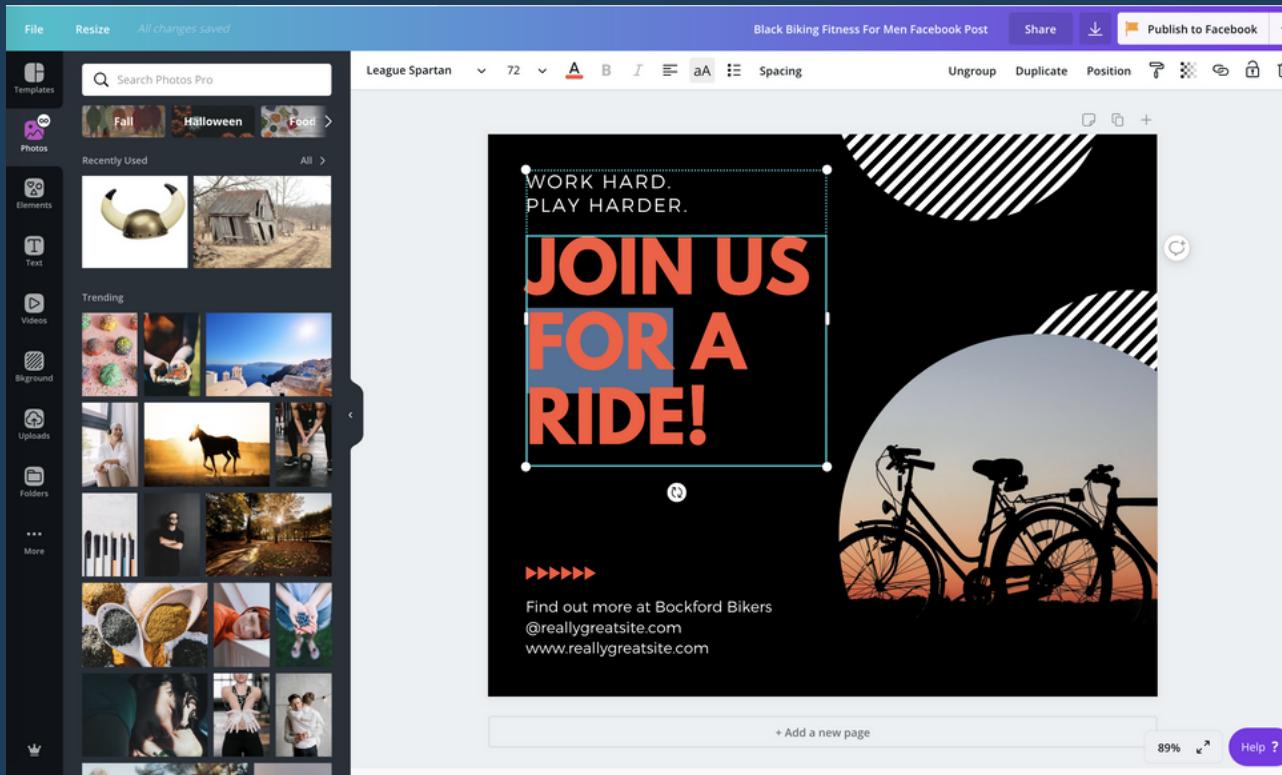


@getwaded

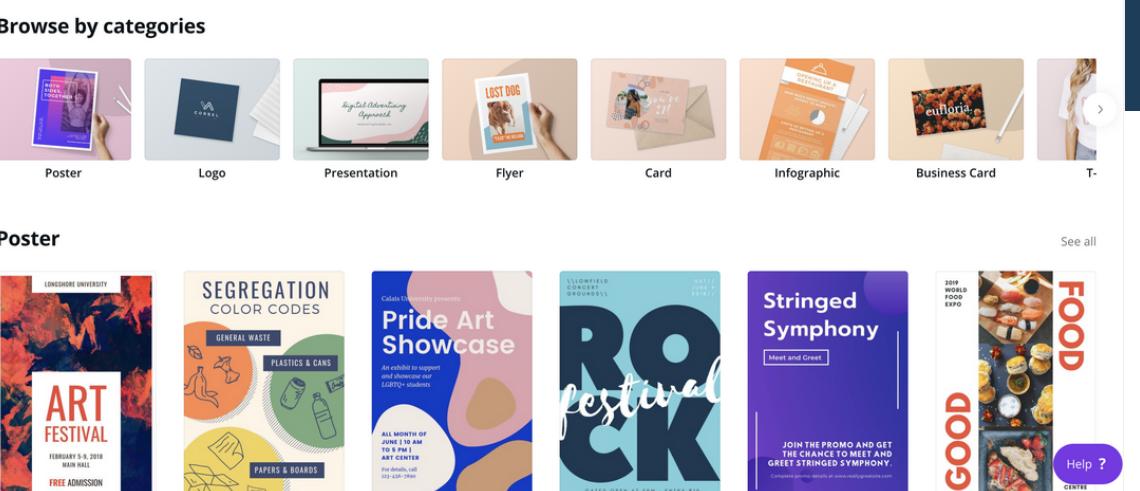
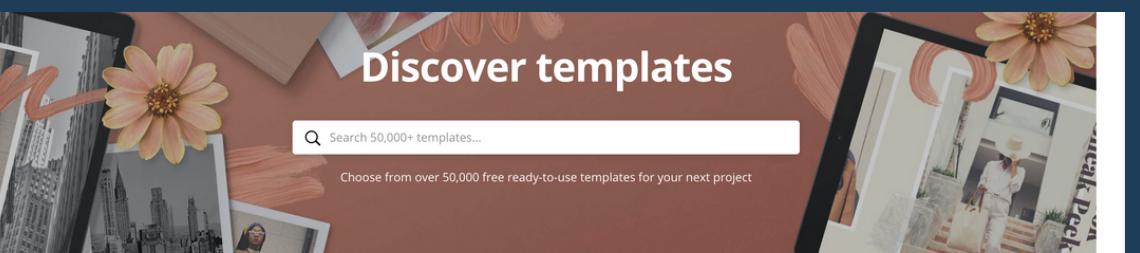


@ArthBaud

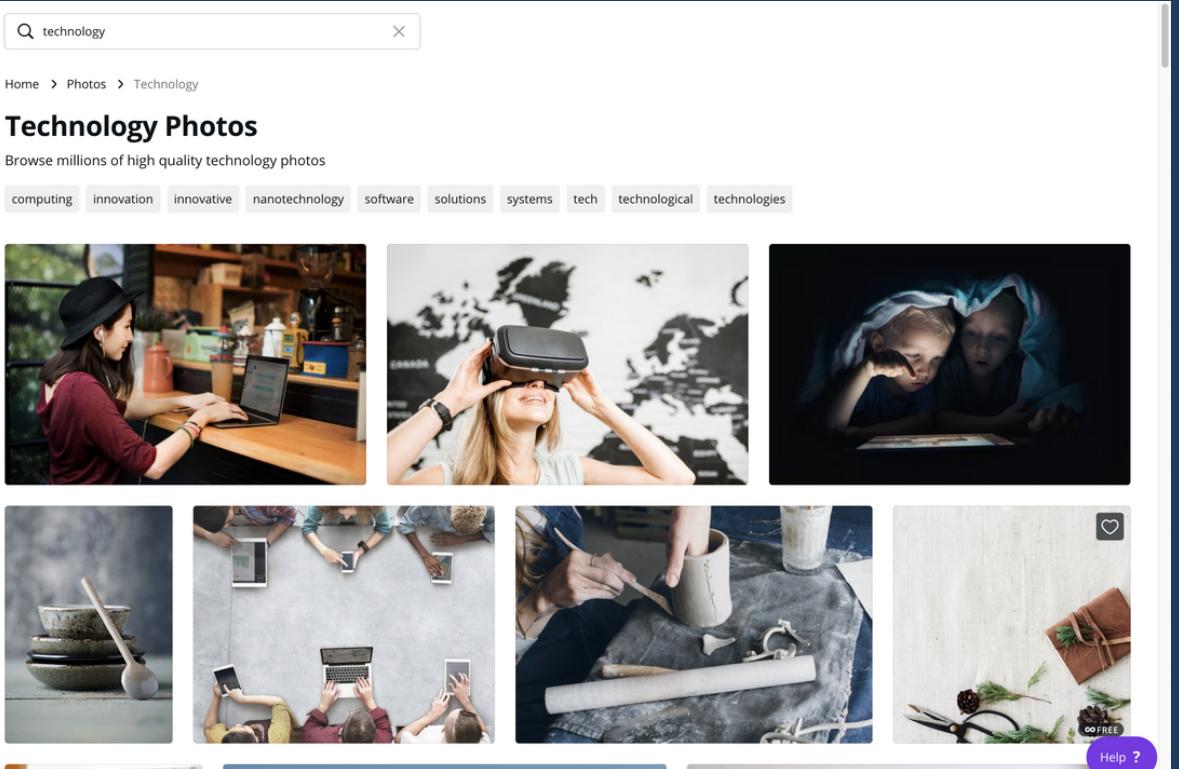
Editor



Templates



Photos

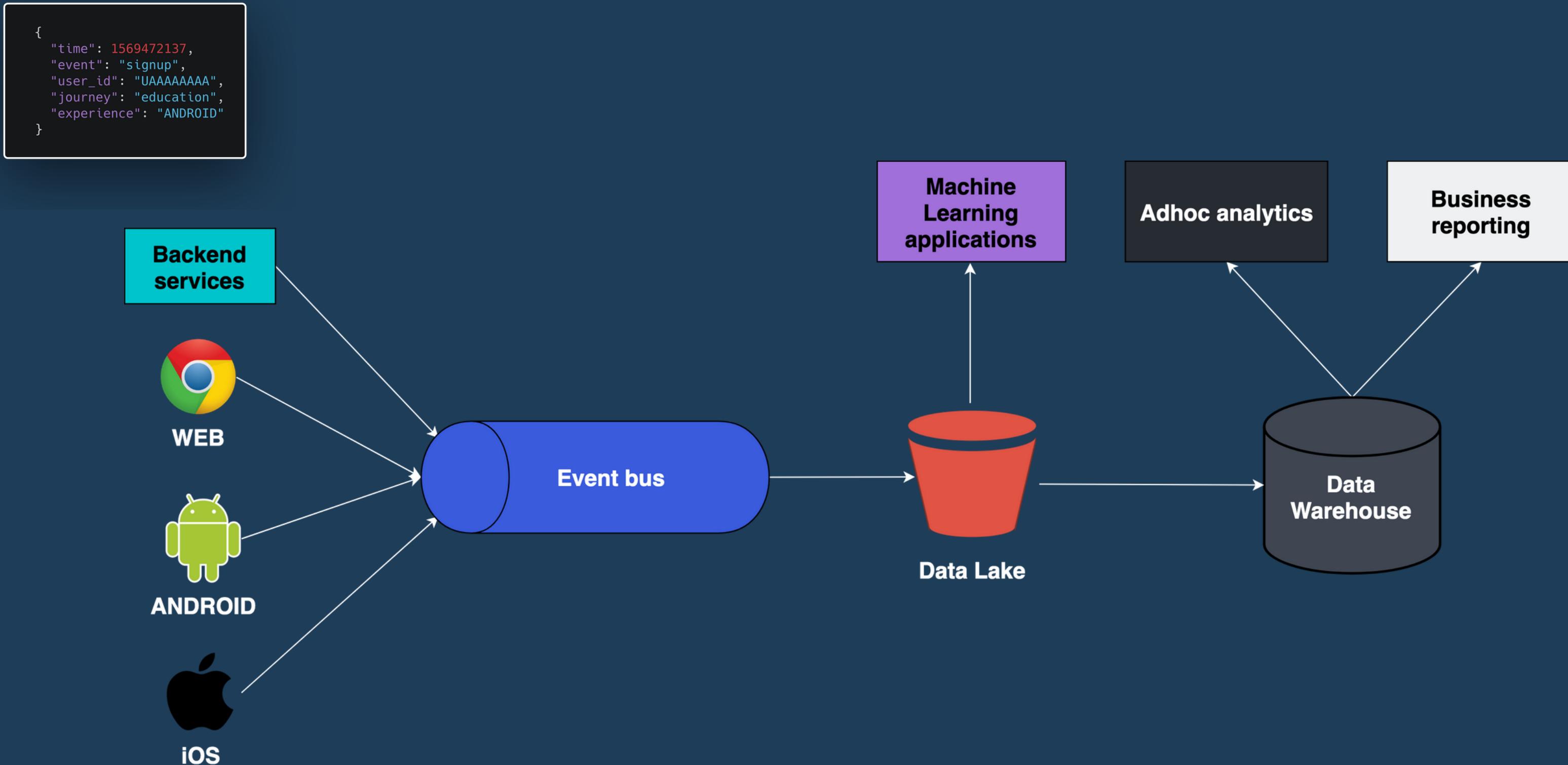


DATA ENG TEAM

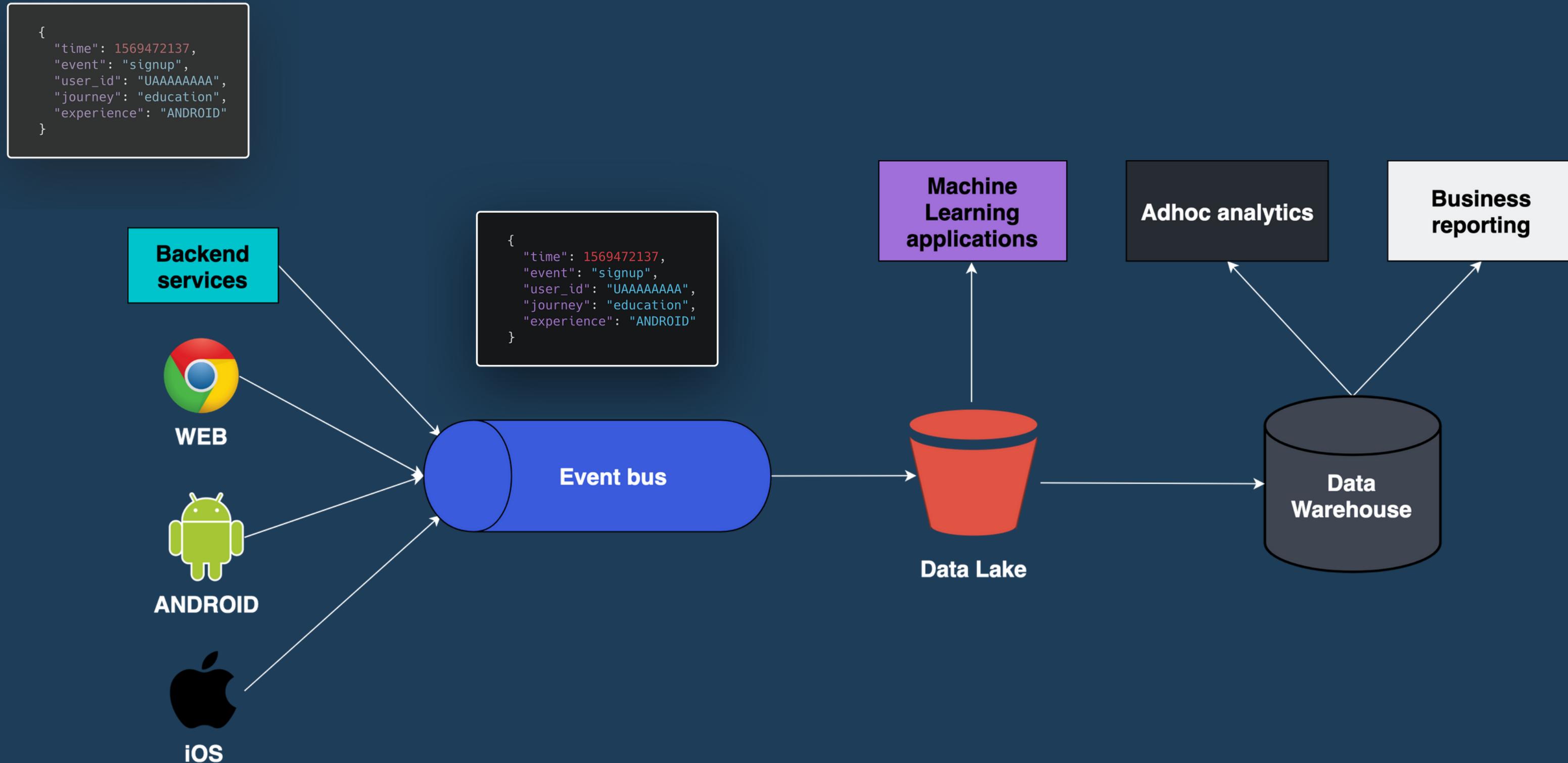
Help Canva make better decisions with data

- Make data easy to produce
- Make data easy to consume
- Data storage and management
- Platform for data pipelines

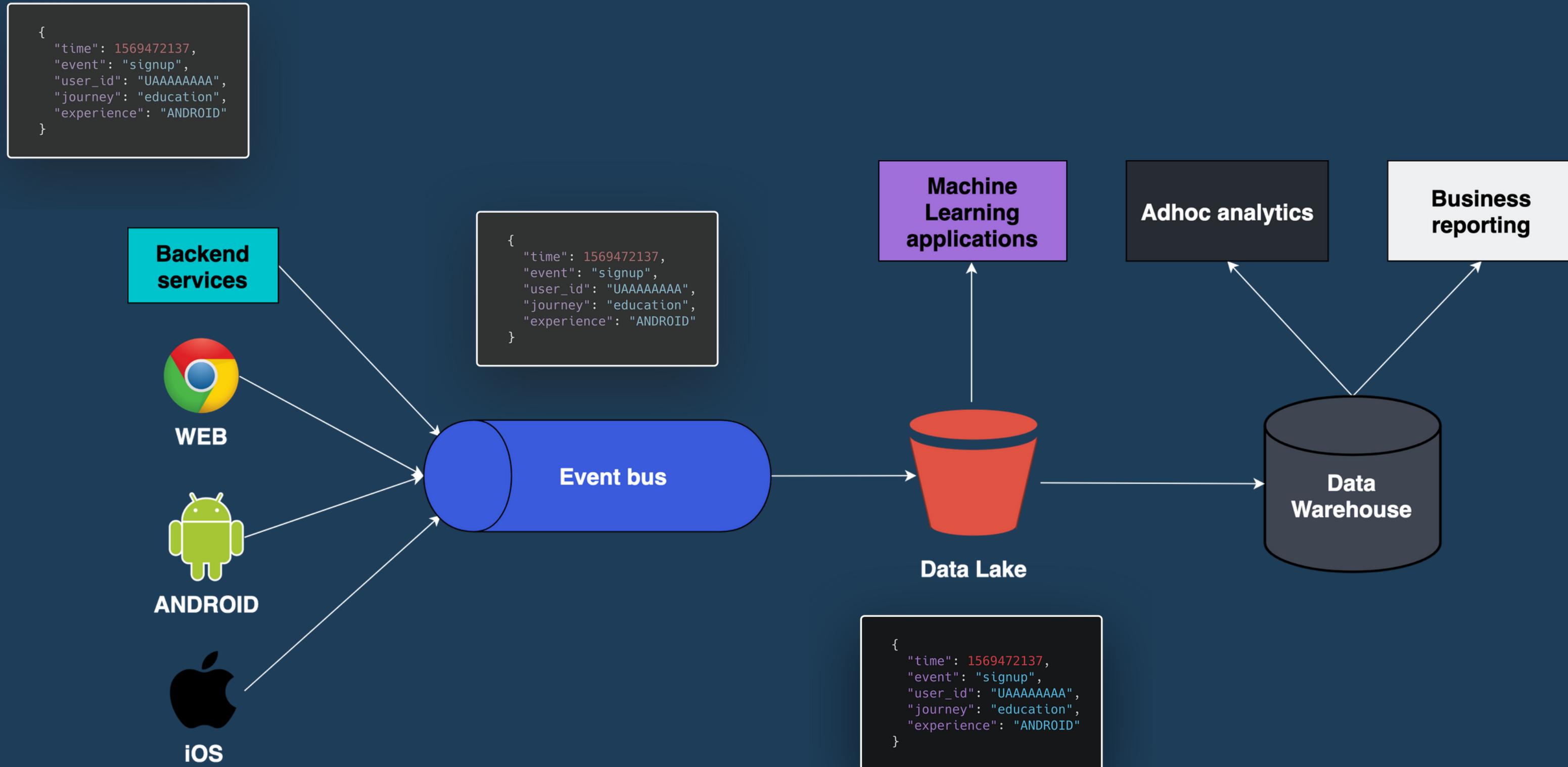
First party event ingestion



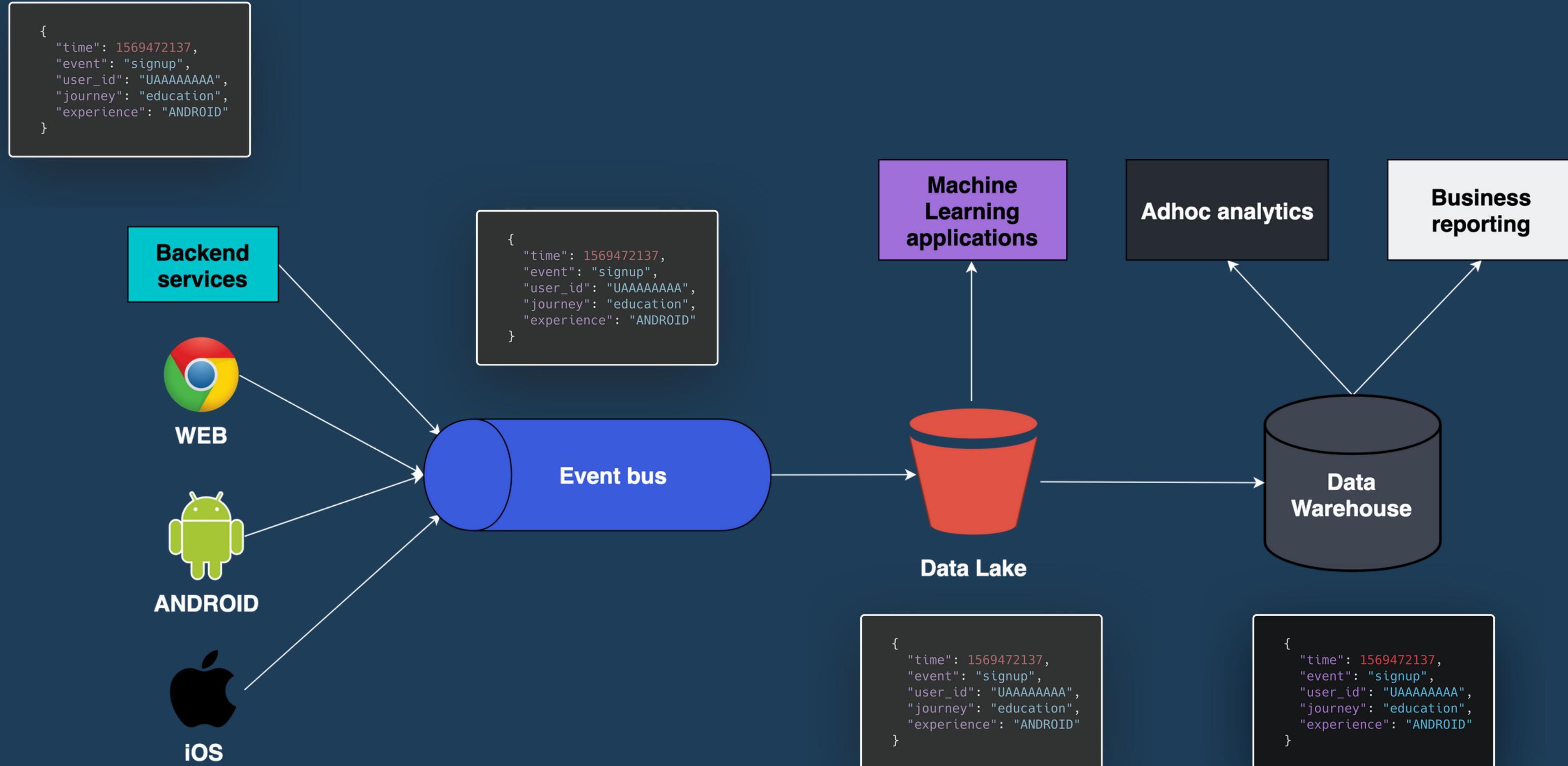
First party event ingestion



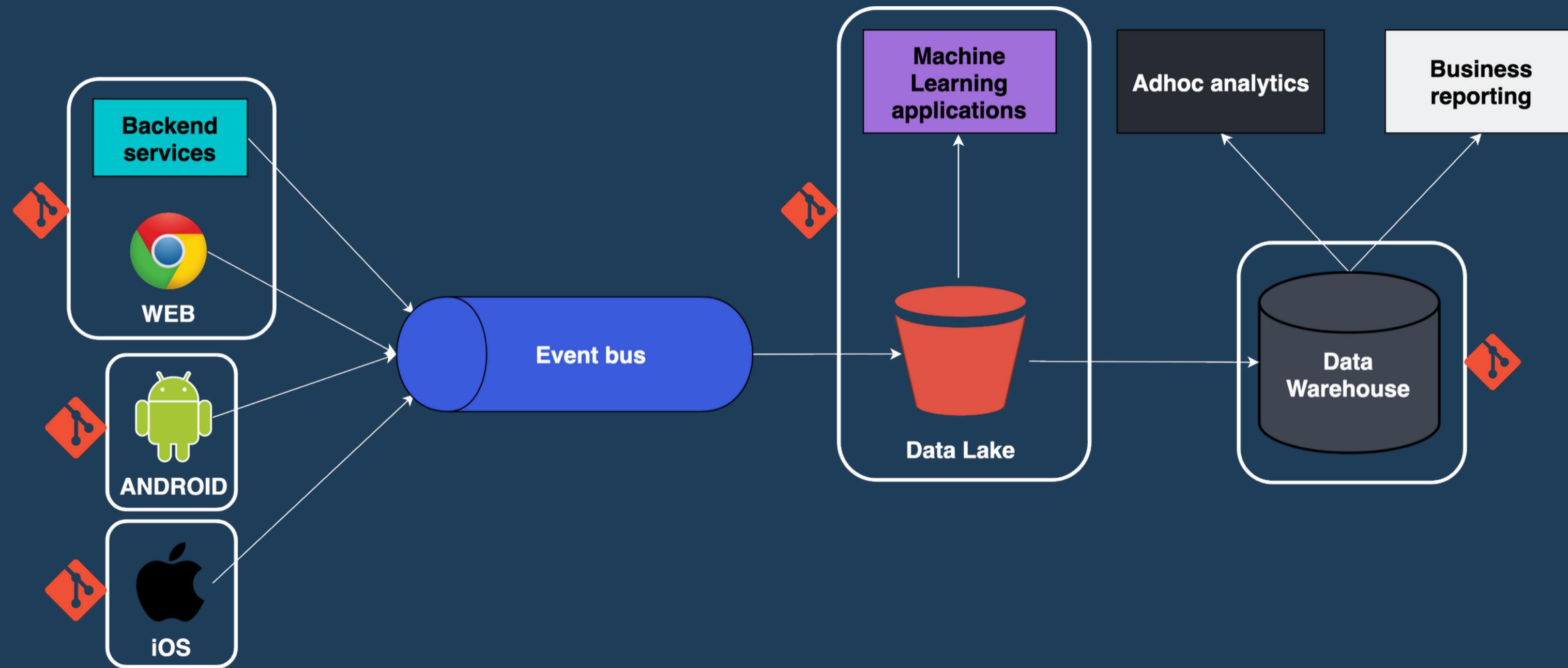
First party event ingestion



First party event ingestion

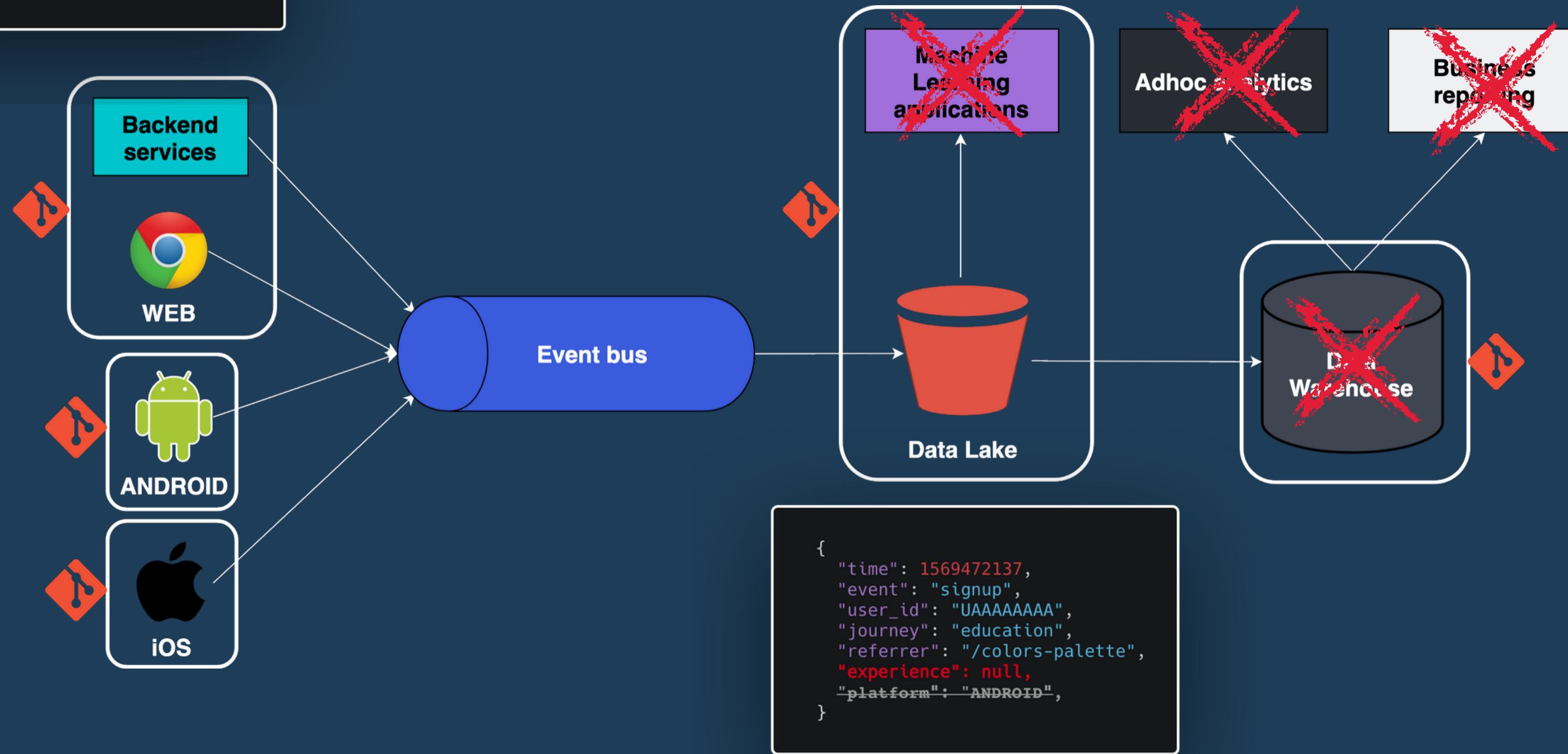


Git repository overlay



Example: Removing a field

```
{  
  "time": 1569472137,  
  "event": "signup",  
  "user_id": "UAAAAAAA",  
  "journey": "education",  
  "referrer": "/colors-palette",  
  - "experience": "ANDROID",  
  + "platform": "ANDROID",  
}
```



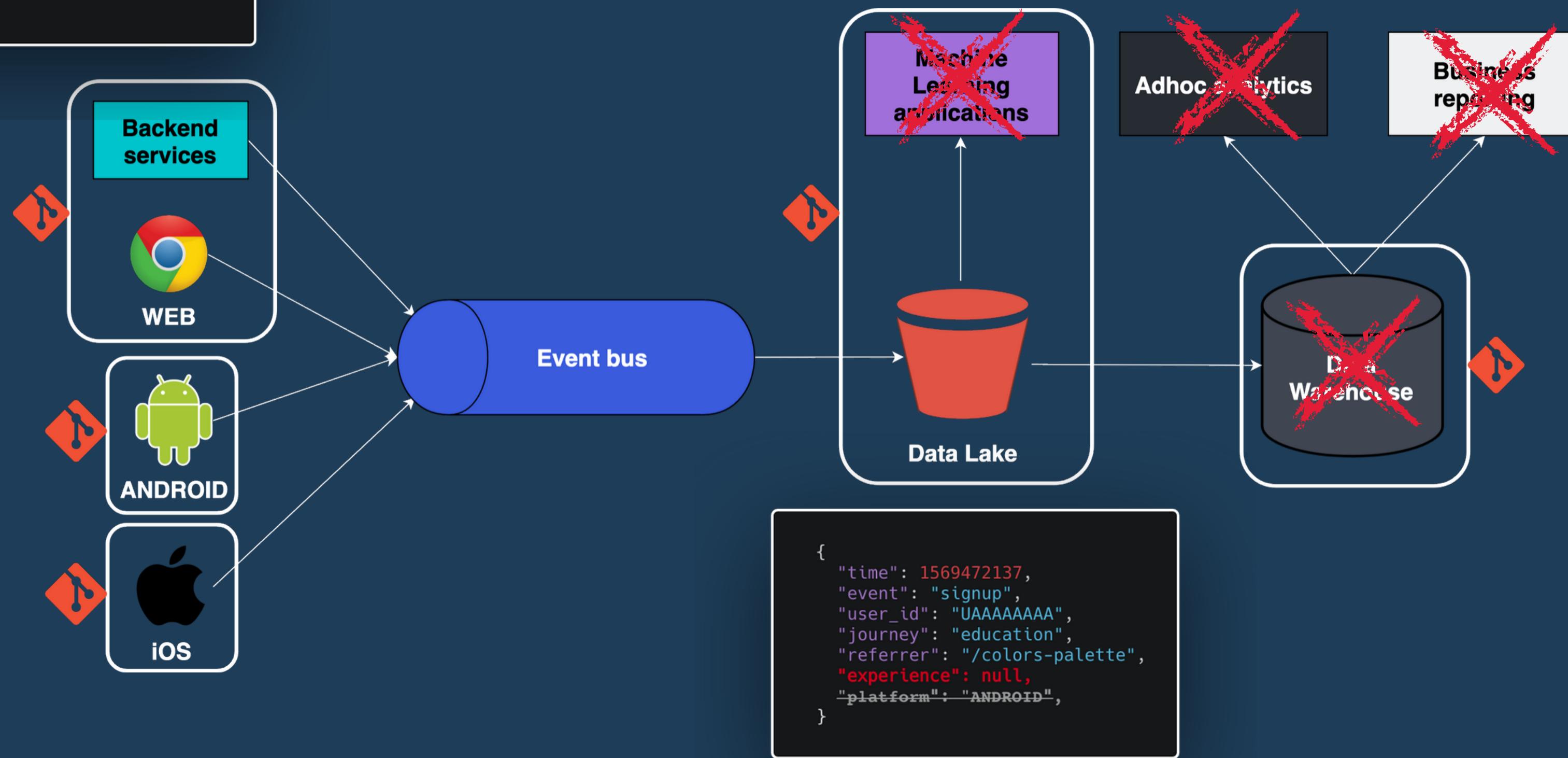
Problems

1. Multiple (conflicting ☹) sources of truth for schemas
2. Schema review and sharing
3. Backwards & forwards compatibility

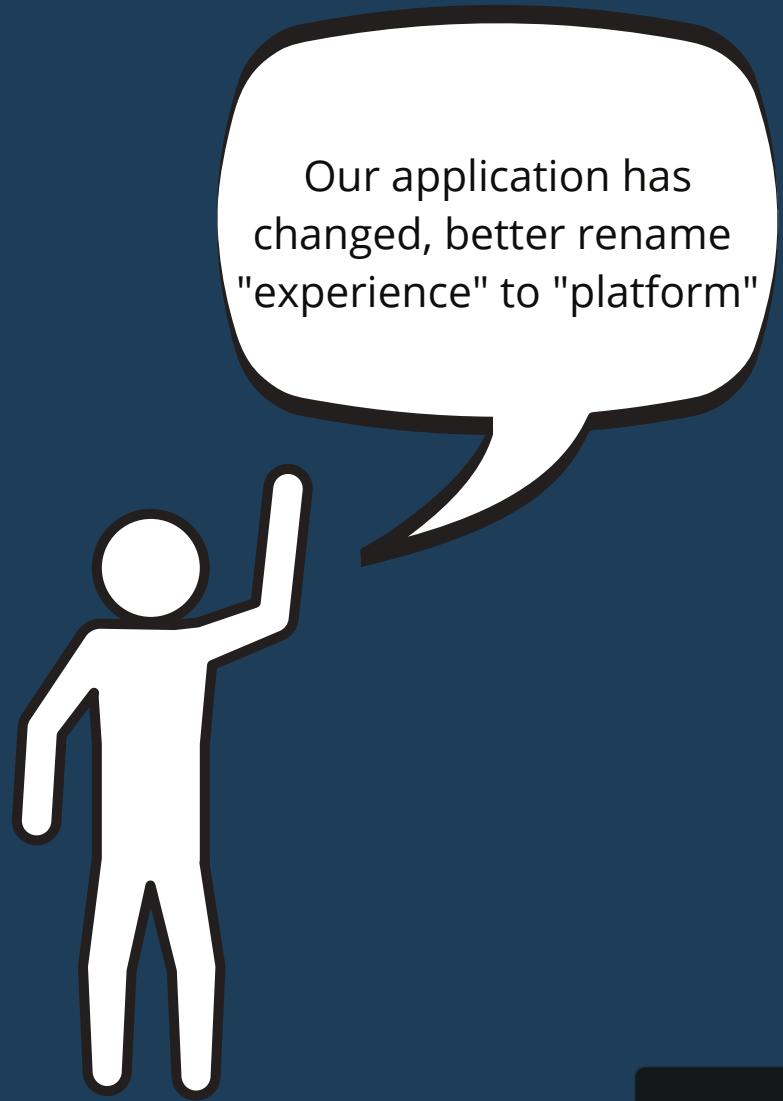
Issue 1 - Multiple sources of truth

Example: Removing a field

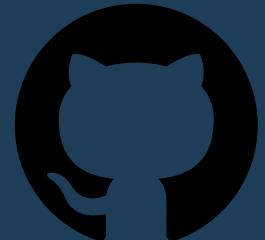
```
{  
  "time": 1569472137,  
  "event": "signup",  
  "user_id": "UAAAAAAA",  
  "journey": "education",  
  "referrer": "/colors-palette",  
  - "experience": "ANDROID",  
  + "platform": "ANDROID",  
}
```



Producer/Consumer



Producer



Repo X/Team A

```
Signup signup = Signup.builder()
    .withTime(1569472137)
    .withEvent("signup")
    .withUserId("AAAAAAAAAA")
    .withJourney("education")
    .withReferrer("/colors-palette")
    - .withExperience("ANDROID")
    + .withPlatform("ANDROID")
    .build();
EventStream output = new EventStream();
signup.writeTo(output);
```

Producer code (Java)

```
{
  "time": 1569472137,
  "event": "signup",
  "user_id": "AAAAAAAAAA",
  "journey": "education",
  "referrer": "/colors-palette",
  - "experience": "ANDROID",
  + "platform": "ANDROID",
}
```

JSON payload

```
@dataclass
class Signup:
    time: int
    event: str
    user_id: str
    journey: str
    referrer: str
    experience: str

input = EventStream()
signup.deserialize(input.read())
```

Consumer code (Python)



Consumer



Repo Y/Team B

Consequences

```
{  
  "time": 1569472137,  
  "event": "signup",  
  "user_id": "UAAAAAAA",  
  "journey": "education",  
  "referrer": "/colors-palette",  
  - "experience": "ANDROID",  
  + "platform": "ANDROID",  
}
```

JSON payload



```
SELECT COUNT(DISTINCT(user_id))  
FROM warehouse.signup  
GROUP BY experience;
```

SQL



```
val signupDF = spark.read  
  .parquet("/event/signup")  
  .as[SignupEvent]  
signupDF  
  .groupBy("experience")  
  .agg(countDistinct("user_id"))
```

Spark



(Silent) Failures

Q: How do we keep production and consumption of events in sync?

A: Schemas

Enters Protobuf

```
syntax = "proto2";

message Signup {
    required int64 time = 1;
    required string event = 2;
    required string user_id = 3;
    required string journey = 4;
    optional string referrer = 5;
    required string experience = 6;
}
```

**Language/platform-neutral extensible
mechanism for serialising structured data**

Interface definition language (IDL)

```
syntax = "proto2";  
  
message Signup {  
    required int64 time = 1;  
    required string event = 2;  
    required string user_id = 3;  
    required string journey = 4;  
    optional string referrer = 5;  
    required string experience = 6;  
}
```

Single schema
definition

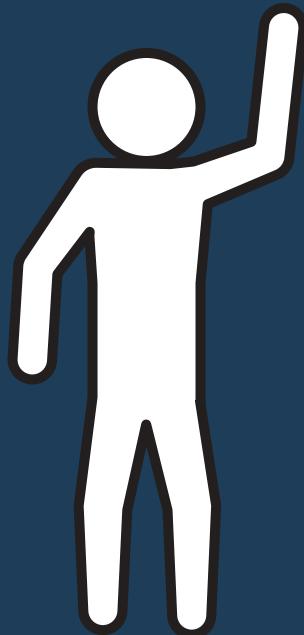


```
Signup signup = Signup.builder()  
    .withTime(1569472137)  
    .withEvent("signup")  
    .withUserId("AAAAAAAA")  
    .withJourney("education")  
    .withReferrer("/colors/palette")  
    .withExperience("ANDROID")  
    .build();
```

```
@dataclass  
class Signup:  
    time: int  
    event: str  
    user_id: str  
    journey: str  
    referrer: str  
    experience: str
```

Multiple code
representation

Serialisation (binary) format



Producer



Repo X

```
Signup signup = Signup.builder()
    .withTime(1569472137)
    .withEvent("signup")
    .withUserId("AAAAAAAAAA")
    .withJourney("education")
    .withReferrer("/colors/palette")
    .withExperience("ANDROID")
    .build();
EventStream output = new EventStream();
signup.writeTo(output)
```

Producer code (Java)

1 0 1 0
0 1 0 1
1 0 1 0



Consumer



Repo Y

```
@dataclass
class Signup:
    time: int
    event: str
    user_id: str
    journey: str
    referrer: str
    experience: str

input = EventStream()
signup.deserialize(input.read())
```

Consumer code (Python)

Vanilla protobuf

- Java
- Python
- Go
- Dart
- C++
- C#
- Binary SerDe

Canva's requirements

- Java
- Swift
- Kotlin
- Scala
- JS/TypeScript
- Json/Jackson library
SerDe

Q: How do you generate custom code from proto files?

A: Build a tool or reuse an existing one

Enters Datumgen

```
syntax = "proto2";

message Signup {
    required int64 time = 1;
    required string event = 2;
    required string user_id = 3;
    required string journey = 4;
    optional string referrer = 5;
    required string experience = 6;
}
```



```
class CaseClassModel {
    public String getName();

    public List<CaseClassField> getFields();
}

class CaseClassField {
    public BooleanisRequired();

    public String getJsonKey();

    public String getName();

    public String getType();
}
```

```
caseClass(caseClass) ::= <<
<if(caseClass.fields)>
case class <caseClass.name>(
    <caseClass.fields:caseClassParamDeclaration(); separator=",\n">
)<endif>
>>

caseClassParamDeclaration(field) ::= <%
<if(field.required)>
@JsonProperty(value = "<field.jsonKey>", required = true)<\ >
<else>
@JsonProperty("<field.jsonKey>")<\ >
<endif>
<field.name>:<\ >
<if(field.required)>
<field.type>
<else>
Option[<field.type>]
<endif>
%>
```

StringTemplate



```
case class Signup(
    @JsonProperty("time") time: Long,
    @JsonProperty("event") event: String,
    @JsonProperty("user_id") userId: String,
    @JsonProperty("journey") journey: String,
    @JsonProperty("referrer") referrer: Option[String],
    @JsonProperty("experience") experience: String,
)
```

Rendering

Improved producer/consumer



Summary - Issue 1

- Multiple definitions of the same structure
 - +Single canonical schema definition
- Handwritten code in 2+ languages
 - +Automated code generation for any language

Issue 2 - Schema review and code sharing

**Having a single schema definition
actually doesn't keep people from
making breaking changes**

Driven by consumers or SME



Review process (denied)



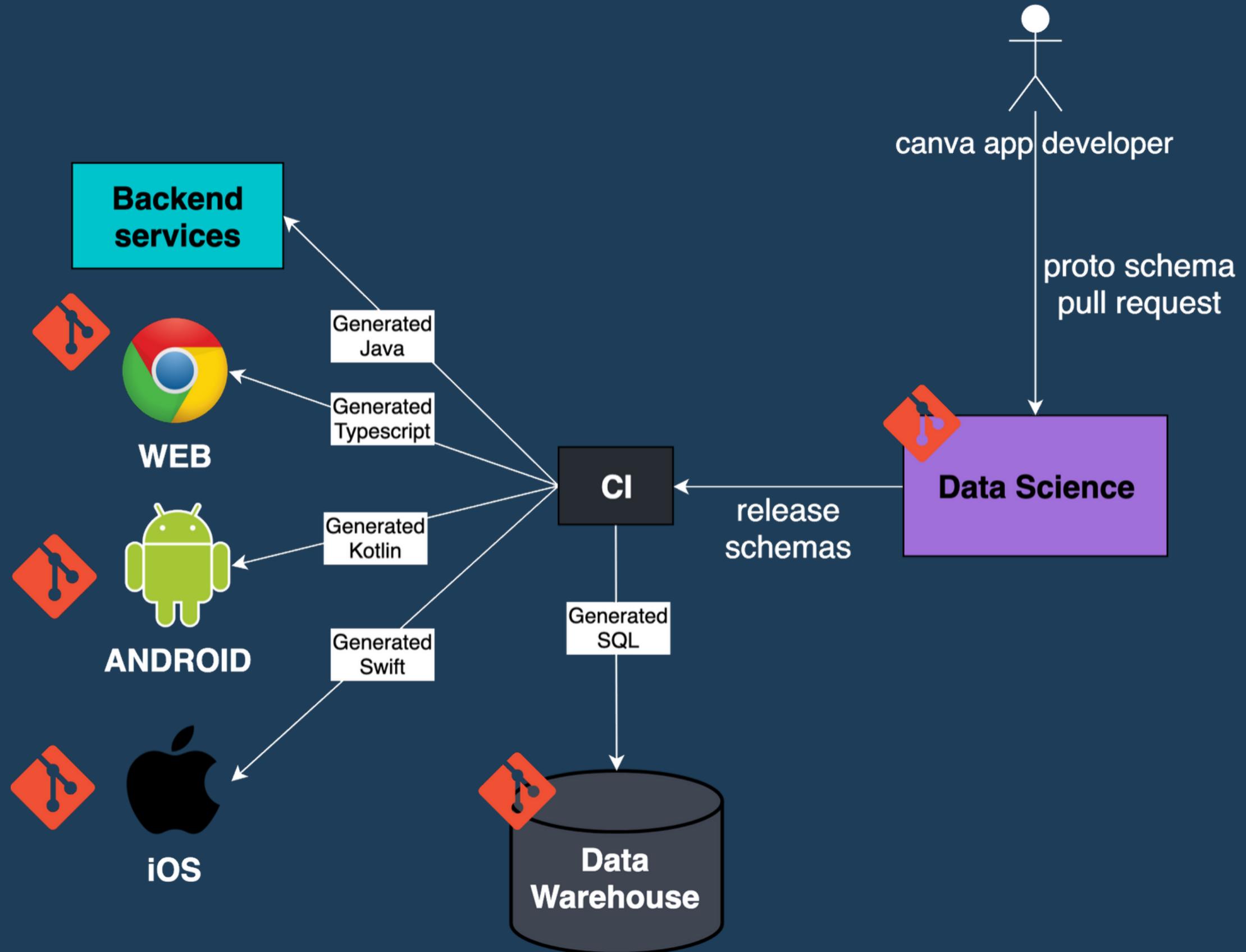
Review process (approved)



**Q: How do you propagate
generated code changes to all its
dependencies**

A: CI

Automated schema sharing



Summary - Issue 2

-Unforeseen/bad schema evolution

+Consumer ownership of schemas, peer review process, guardians of compatibility

-Generated code out of sync

+Automated code propagation

Issue 3 - Backwards and forwards compatibility

Maintaining full transitive compatibility

Compatibility types

Compatibility Type	Changes allowed	Check against which schemas	Upgrade first
BACKWARD	<ul style="list-style-type: none">• Delete fields• Add optional fields	Last version	Consumers
BACKWARD_TRANSITIVE	<ul style="list-style-type: none">• Delete fields• Add optional fields	All previous versions	Consumers
FORWARD	<ul style="list-style-type: none">• Add fields• Delete optional fields	Last version	Producers
FORWARD_TRANSITIVE	<ul style="list-style-type: none">• Add fields• Delete optional fields	All previous versions	Producers
FULL	<ul style="list-style-type: none">• Modify optional fields	Last version	Any order
FULL_TRANSITIVE	<ul style="list-style-type: none">• Modify optional fields	All previous versions	Any order
NONE	<ul style="list-style-type: none">• All changes are accepted	Compatibility checking disabled	Depends

BACKWARD COMPATIBILITY

New consumers need to be able to read both new messages and the previous message version.

Makes sense for:

- RPC services
- Regular release schedule
- Monorepos

Compatibility Type	Changes allowed	Check against which schemas	Upgrade first
BACKWARD	<ul style="list-style-type: none">• Delete fields• Add optional fields	Last version	Consumers
BACKWARD_TRANSITIVE	<ul style="list-style-type: none">• Delete fields• Add optional fields	All previous versions	Consumers
FORWARD	<ul style="list-style-type: none">• Add fields• Delete optional fields	Last version	Producers
FORWARD_TRANSITIVE	<ul style="list-style-type: none">• Add fields• Delete optional fields	All previous versions	Producers
FULL	<ul style="list-style-type: none">• Modify optional fields	Last version	Any order
FULL_TRANSITIVE	<ul style="list-style-type: none">• Modify optional fields	All previous versions	Any order
NONE	<ul style="list-style-type: none">• All changes are accepted	Compatibility checking disabled	Depends

FORWARD

New messages need to be able to be read by **old** consumers.

Makes sense for:

- Decoupled services reading from an event bus (eg. Kafka)
- Regular release schedule
- Monorepos

Compatibility Type	Changes allowed	Check against which schemas	Upgrade first
BACKWARD	<ul style="list-style-type: none">• Delete fields• Add optional fields	Last version	Consumers
BACKWARD_TRANSITIVE	<ul style="list-style-type: none">• Delete fields• Add optional fields	All previous versions	Consumers
FORWARD	<ul style="list-style-type: none">• Add fields• Delete optional fields	Last version	Producers
FORWARD_TRANSITIVE	<ul style="list-style-type: none">• Add fields• Delete optional fields	All previous versions	Producers
FULL	<ul style="list-style-type: none">• Modify optional fields	Last version	Any order
FULL_TRANSITIVE	<ul style="list-style-type: none">• Modify optional fields	All previous versions	Any order
NONE	<ul style="list-style-type: none">• All changes are accepted	Compatibility checking disabled	Depends

<https://docs.confluent.io/current/schema-registry/avro.html#compatibility-types>

Full transitive compatibility

What:

Consumers need to be able to read:

- the current message version,
- all the previous versions,
- all future versions.

Why:

Need to be able to query and make comparisons against long term historical data

How:

- Only add optional fields, or remove unused optional fields
- **No enums!**

CAVEATS*

You can cheat these compatibility rules under 3 circumstances:

1. You backfill all old data into the new schema and migrate all downstream users
2. You are can tolerate data loss when a new producer or consumer is deployed
3. You can guarantee you have no downstream users (yet)

**How do we maintain full
transitive compatibility?**

Protolock

github.com/nilslice/protolock

"Protocol Buffer companion tool. Track your .proto files and prevent changes to messages and services which impact API compatibility."

signup_event.proto

```
syntax = "proto2";

message Signup {
    required int64 time = 1;
    required string event = 2;
    required string user_id = 3;
    required string journey = 4;
    optional string referrer = 5;
    required string experience = 6;
}
```

> protolock init

proto.lock

```
{
  "definitions": [
    {
      "protopath": "signup_event.proto",
      "def": {
        "messages": [
          {
            "name": "Signup",
            "fields": [
              {
                "id": 1,
                "name": "time",
                "type": "int64"
              },
              {
                "id": 2,
                "name": "event",
                "type": "string"
              },
              ...
              ...
              ...
              {
                "id": 6,
                "name": "experience",
                "type": "string"
              }
            ]
          }
        ],
        "package": {
          "name": "com.canva.example"
        }
      }
    }
  ]
}
```

signup_event.proto

```
syntax = "proto2";

message Signup {
    required int64 time = 1;
    required string event = 2;
    required string user_id = 3;
    required string journey = 4;
    optional string referrer = 5;
    -required string experience = 6;
    +optional string platform = 7;
}
```

> protolock commit



CONFLICT: "Signup" ID: "6"
has been removed
[signup_event.proto]

CONFLICT: "Signup" field:
"experience" has been
removed
[signup_event.proto]

Protolock Compatibility Rules

- No changing field IDs
- No changing field types
- No changing field names
- No removing required fields
- No adding new required fields

Social engineering via red bubbles

rename user city location field

Build #13934 | wade-breaking-schema-change | ⚡ 31792f796e (Pull Request #2885) Run

 Verify Regen Twirl  Verify DS Schemas  Verify Schema Lock  Datumgen fixtures

Bonus benefits of generating code from schemas

- Generated producers and consumers
- Standardise events across apps / platforms:
eg. web, backend, iOS, Android
- Abstract the event bus
- Abstract the storage layer
- Abstract monitoring
- Wide ranging refactors become trivial
- Better schema design when reviews are centralised
- Easy to assign data owners

Cons of generating code from schemas

- Very hard to retrofit into an existing system
- Requires buy-in from all data producers
- Schema sharing across multiple repos is eventually consistent
- Big tooling investment

Summary

Before:

- Multiple manually written sources of truth
- No best practice around event design
- Breaking changes everywhere

Summary

Before:

- Multiple manually written sources of truth
- No best practice around event design
- Breaking changes everywhere

After:

- Generated producer and consumer clients from proto
- Centralised schema review process and publishing mechanism
- Automated CI checks enforce full transitive compatibility

We're hiring



- Data Engineers
- Machine Learning Engineers
- Data Analysts

<https://www.canva.com/careers/jobs/>

Questions