# Data Pipeline Design

*Advanced Python Programming*

—

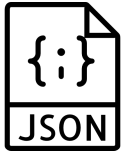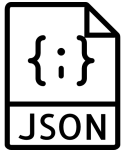Valerio Velardo - The Sound of AI

# What do we want to build?

A data pipeline that reads product records from json files, processes them, and stores the processed data in a DB.
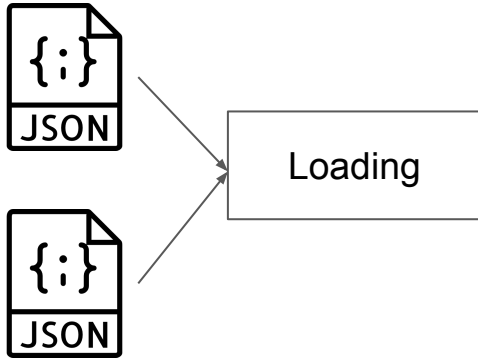
# Example product json file

```json
{
    "name": "product_1",
    "currency": "dollar",
    "price": 100
}
```
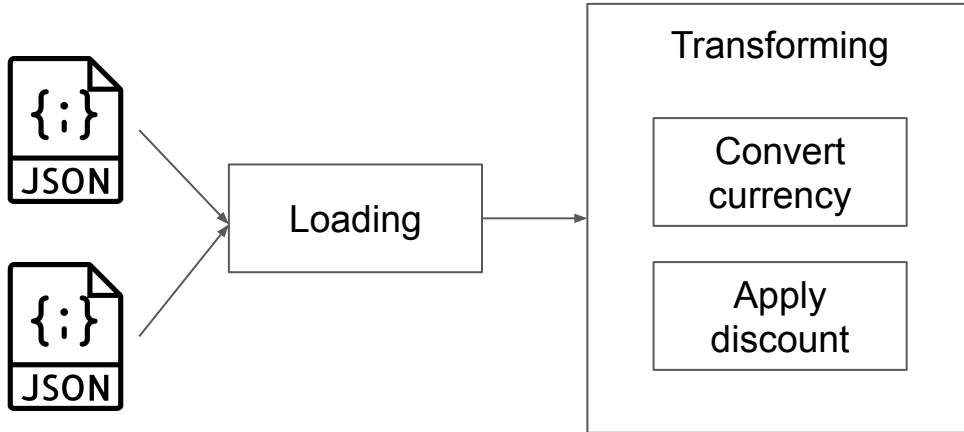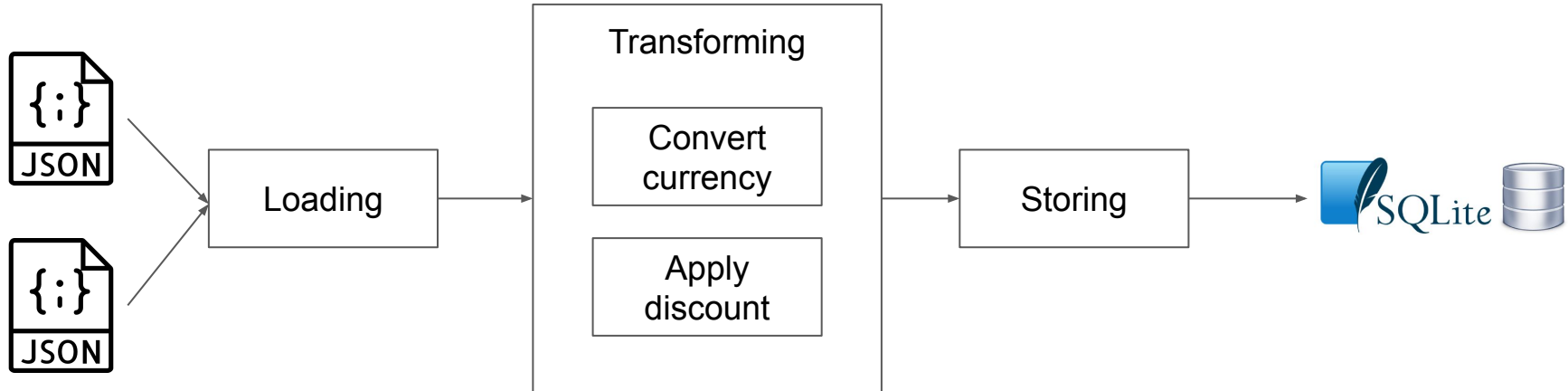
# Programme workflow

# Programme workflow
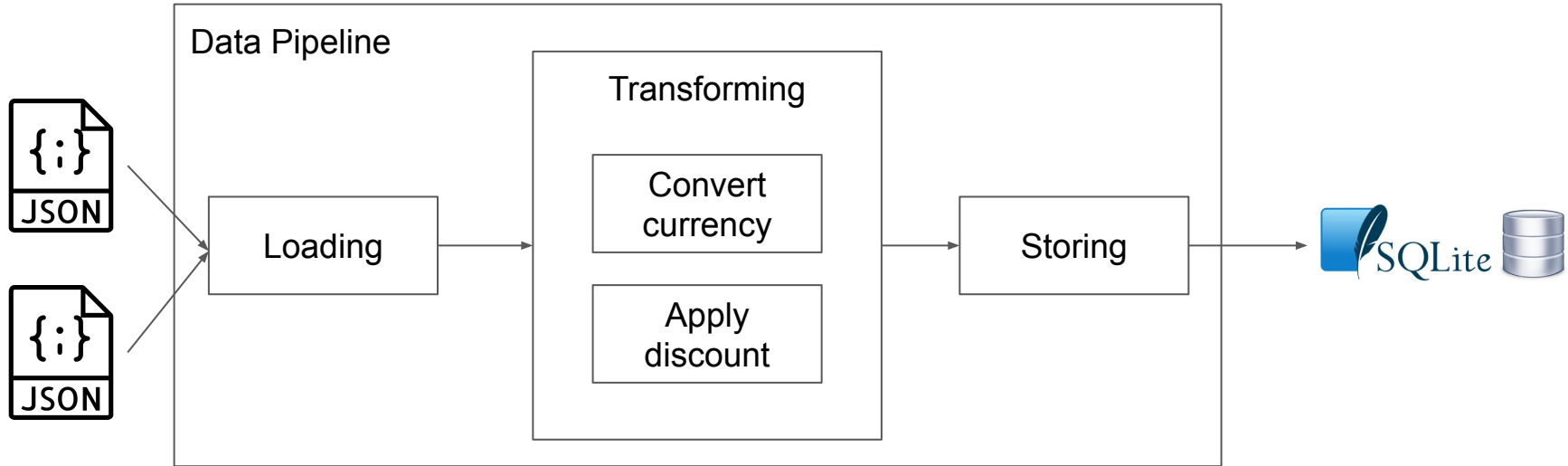
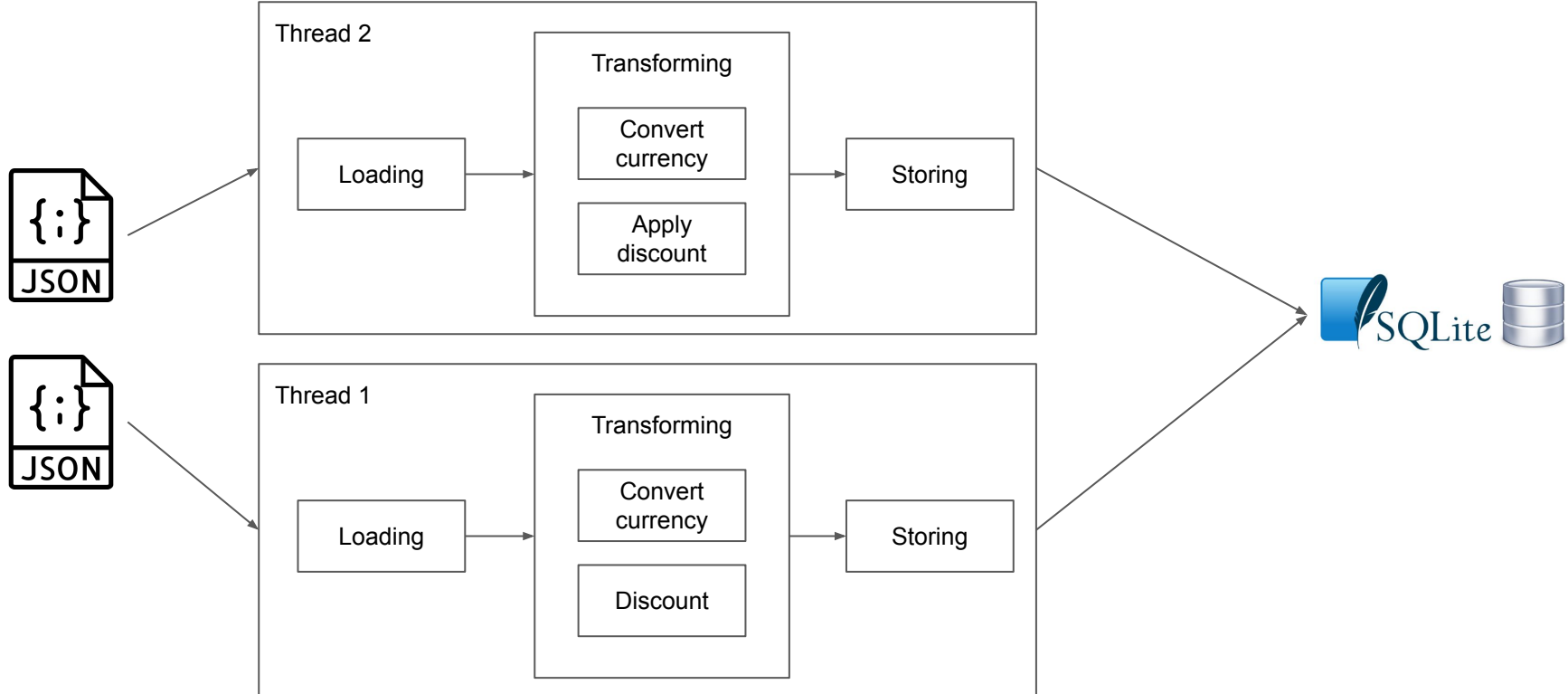# Programme workflow

# Programme workflow

# Programme workflow

# Programme workflow

# Programme workflow with threading

# Product DTO

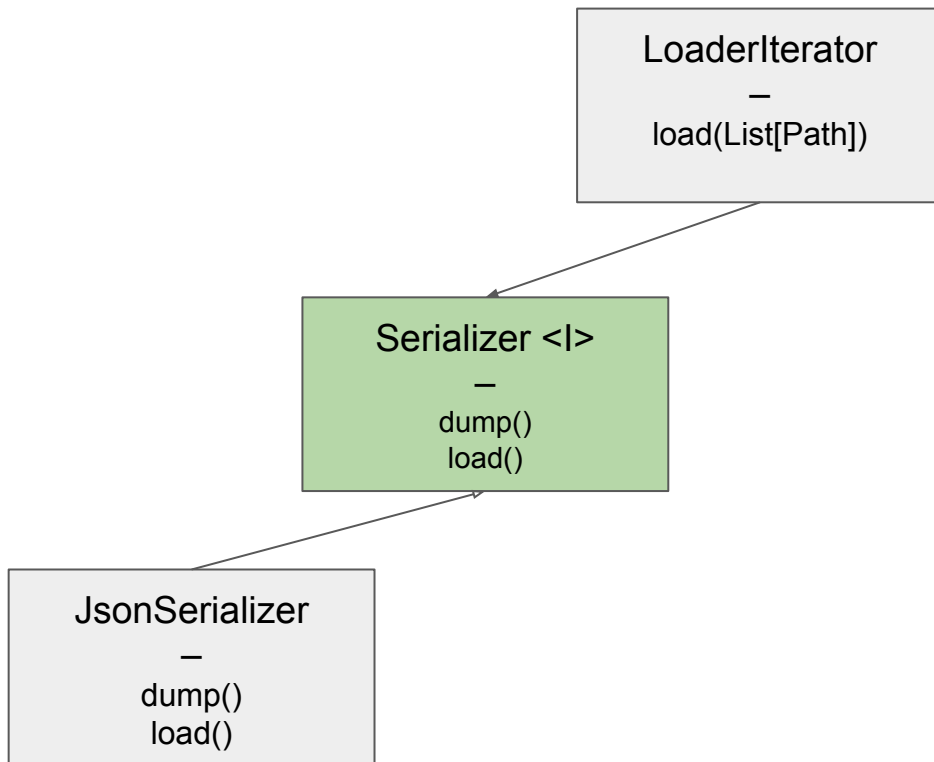Product(pydanctic.BaseModel)
—
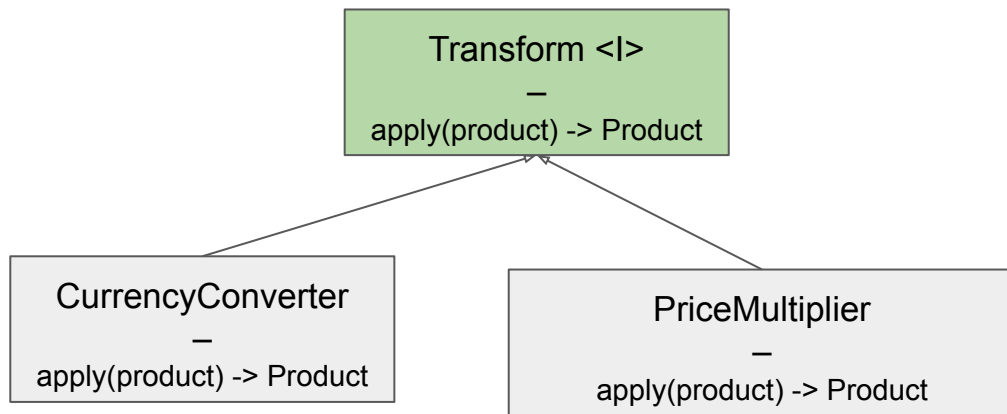name
currency
price

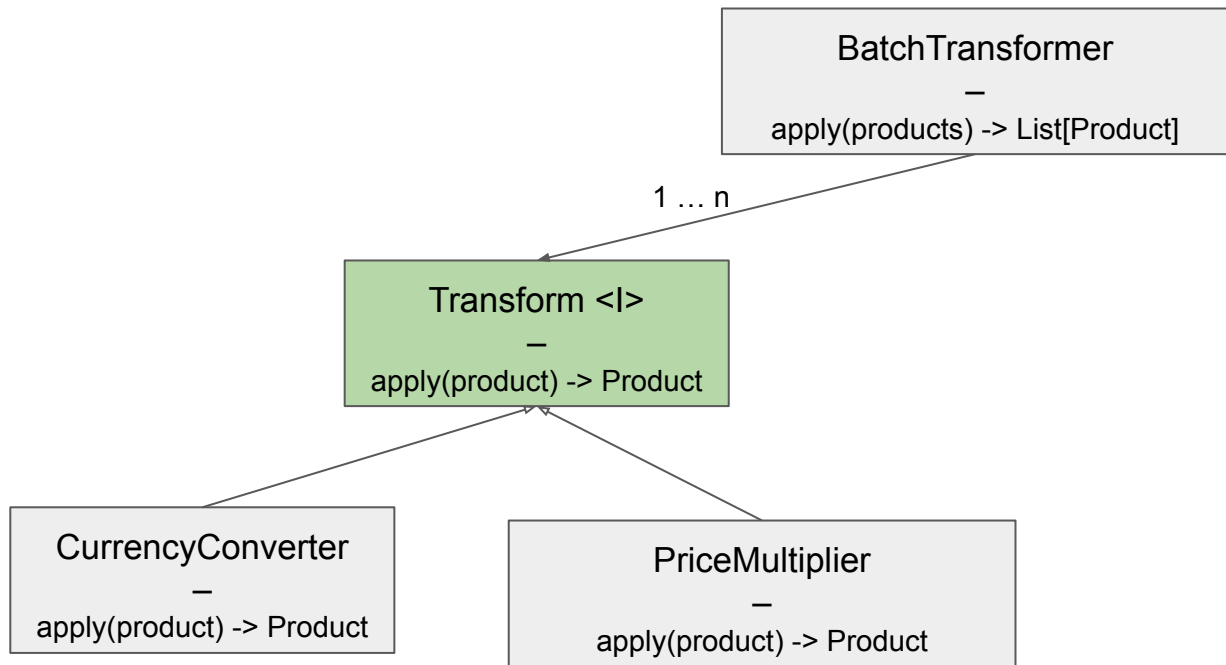# Loading component

LoaderIterator
–
load(List[Path])

# Loading component

# Transforming component

```
          Transform <I>
                –
       apply(product) -> Product
```

```
  CurrencyConverter              PriceMultiplier
         –                              –
apply(product) -> Product    apply(product) -> Product
```

# Transforming component

# Storing component

SqLiteProductStorer
–
store(sqlite.Cursor,
List[products])

SQLiteContextManager

# Data pipeline component

# Threading component

# What's next?

1. Product DTO

2. Loading component

3. Transforming component

4. Storing component

5. Data pipeline

6. Multi threading

7. Running the programme