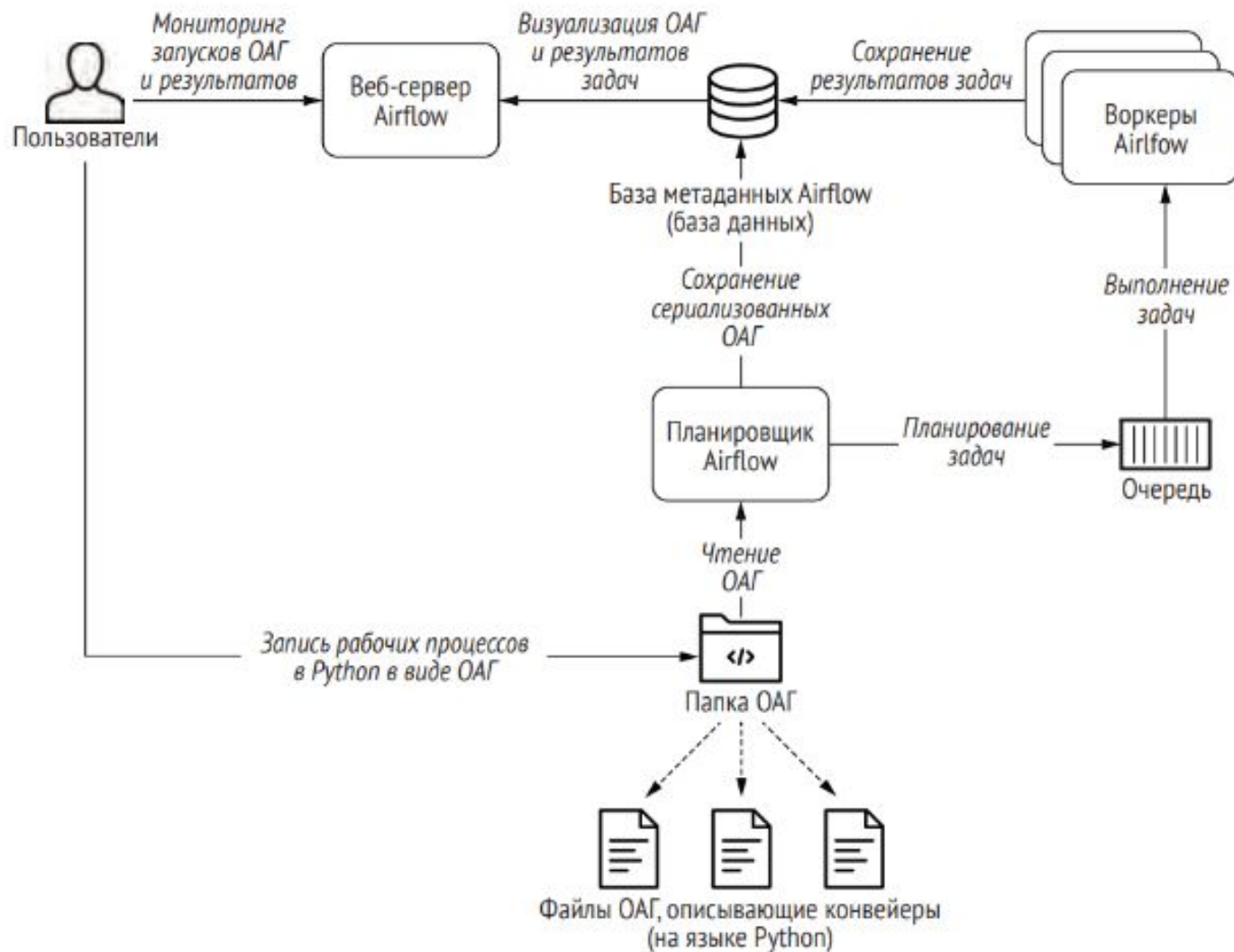

Обработка больших данных с Apache Airflow (ч. 2)

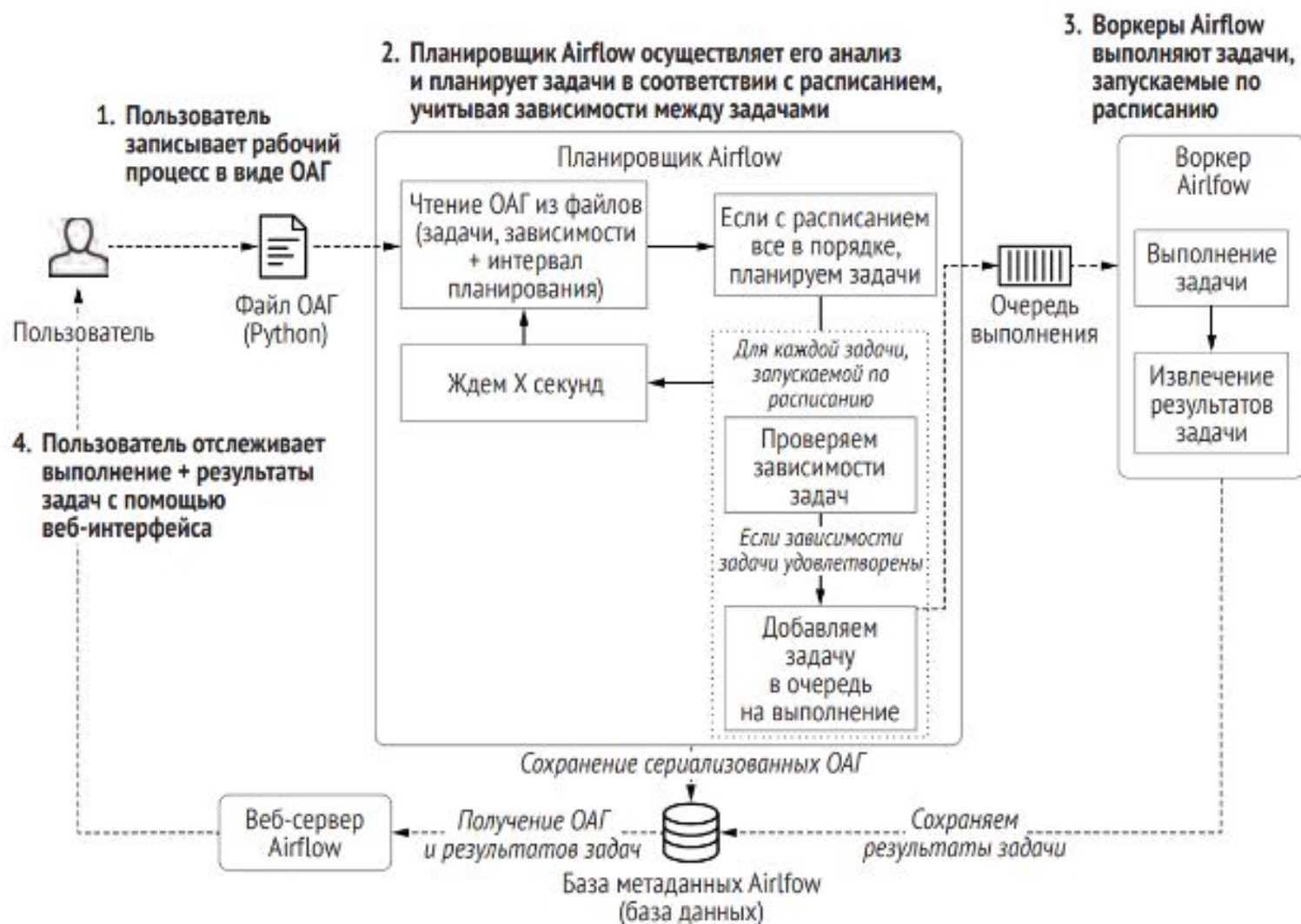
Астахов Сергей ИУ6-22М

Agenda

- Планирование и выполнение конвейеров
- Инкрементальная загрузка и обратное заполнение
- Атомарность и идемпотентность задач
- Ветвление

Планирование и выполнение конвейеров





Инкрементальная загрузка и обратное заполнение

Catchup = true (default)

AirFlow начинает обработку, включая прошлые интервалы
с *backfill*

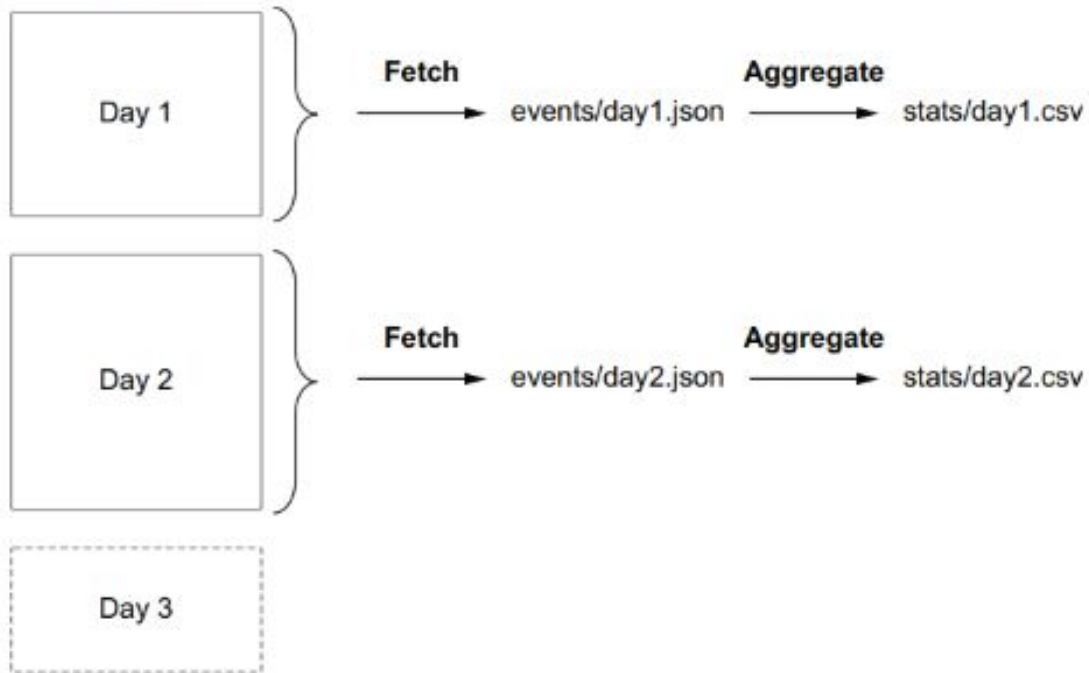


Catchup = false

AirFlow начинает обработку
с текущего интервала



Events



Атомарность и идемпотентность

Не атомарно

```
def _calculate_stats(**context):
    """Calculates event statistics."""
    input_path = context["templates_dict"]["input_path"]
    output_path = context["templates_dict"]["output_path"]

    events = pd.read_json(input_path)
    stats = events.groupby(["date",
                           "user"]).size().reset_index()
    stats.to_csv(output_path, index=False)

    send_stats(stats, email="user@example.com")
```

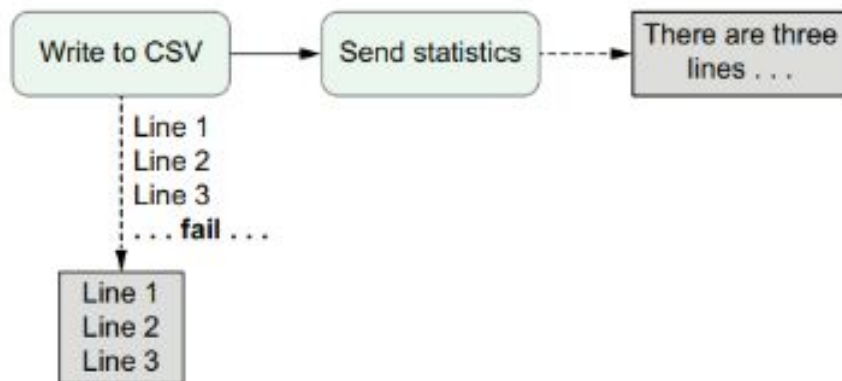
Атомарно

```
def _send_stats(email, **context):
    stats =
    pd.read_csv(context["templates_dict"]["stats_path"]
    ")
    email_stats(stats, email=email)

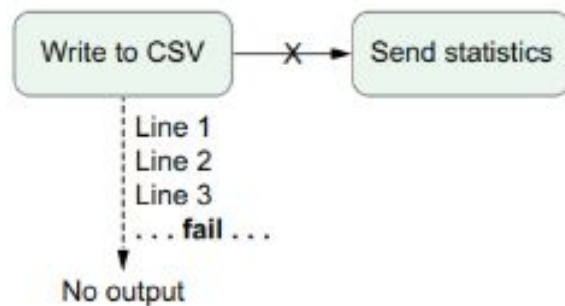
send_stats = PythonOperator(
    task_id="send_stats",
    python_callable=_send_stats,
    op_kwargs={"email": "user@example.com"},
    templates_dict={"stats_path":
    "data/stats/{[ds]}.csv"},
    provide_context=True,
    dag=dag,
)
```

calculate_stats >> send_stats

Non-atomic operation

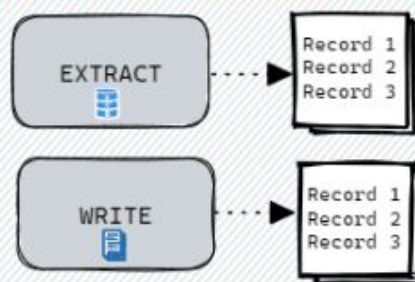


Atomic operation

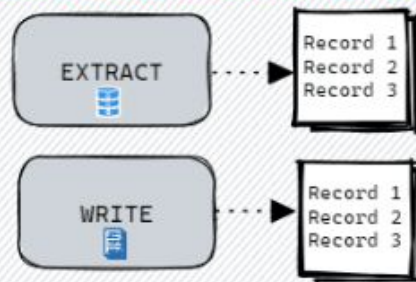


RUN #1

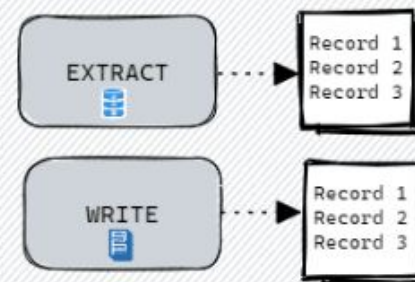
Идемпотентная задача, каждый раз выдающая один и тот же результат



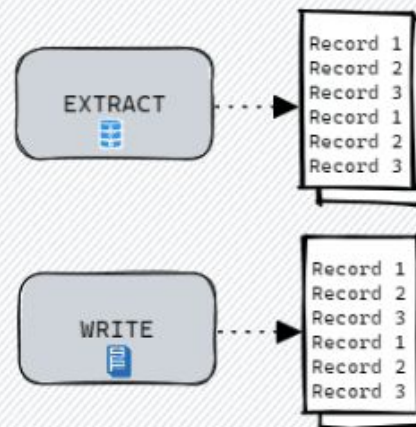
RUN #2



RUN #1



RUN #2



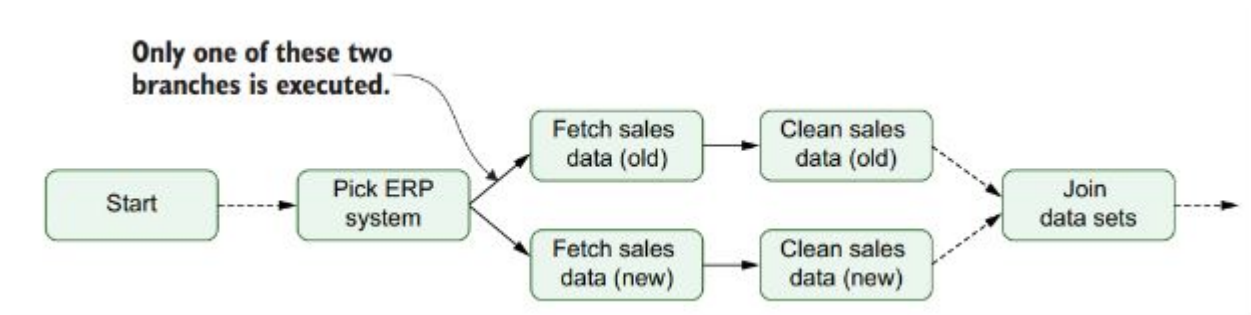
Неидемпотентная задача,
выдающая повторы в результате

IDEMPOTENCE

WHEN PERFORMING AN OPERATION AGAIN GIVES THE SAME RESULT

HTTP METHOD	IDEMPOTENCE	SAFETY
GET	YES	YES
HEAD	YES	YES
PUT	YES	NO
DELETE	YES	NO
POST	NO	NO
PATCH	NO	NO

Ветвление



```
def _pick_erp_system(**context):  
    if context["execution_date"] < ERP_SWITCH_DATE:  
        return "fetch_sales_old"  
    else:  
        return "fetch_sales_new"
```

```
sales_branch = BranchPythonOperator(  
    task_id='sales_branch',  
    provide_context=True,  
    python_callable=_pick_erp_system,  
)
```

```
sales_branch >> [fetch_sales_old, fetch_sales_new]
```

- all_success
- all_failed
- all_done
- all_skipped
- one_failed
- one_success
- one_done
- none_failed
- none_failed_min_one_success
- none_skipped
- always

- BranchPythonOperator
 - ShortCircuitOperator
 - BranchSQLOperator
 - BranchDayOfWeekOperator
 - BranchDateTimeOperator
-

**Спасибо
за внимание**



**Спасибо
за внимание**