

# AMPEL

- 1. Motivation
- 2. Concepts
- 3. Selected tools
- 4. Lessons Learned



Jakob (Nordin alt. van Santen)



# I would like...

| High-thro     | oughput            |                       | R          | Repeatability | F  | AIR                |
|---------------|--------------------|-----------------------|------------|---------------|----|--------------------|
| Real-time     | Large sam          | ple                   | <b>TTT</b> | - 1:15        | Ha | ave only one       |
| Fast response | e More p<br>more d | oixels,<br>lishes     | w ny       | Test statisti | .C | Green<br>computing |
|               | Flexibility        | Software<br>developme | ent        | Statistics    |    |                    |
|               | Big Data           |                       | Mach       | nine Learning | 5  |                    |





Jakob (Nordin alt. van Santen)



Jakob (Nordin alt. van Santen)



Zeuthen Data Science Seminar

Jakob (Nordin alt. van Santen)



<u>AMPEL</u> is a *modular* and *scalable* platform with explicit *provenance* tracking, suited for systematically processing large (possibly complex and heterogeneous) data streams.

This includes selecting, analyzing, updating, combining, enriching and reacting to data.

**Code-To-Data**: Teams develop analysis algorithms which are hosted in AMPEL and exposed to data streams in a high performance computer center.





### Sample use case: Neutrino-TDE coincidence searches



### Match TDE model to real-time stream. Cross-match to IceCube alerts.



Jakob (Nordin alt. van Santen)

### Selected conceptual components

# 1 Deconstructing scientific workflows.

Modularity requires rules for how modules interact.

Theory: workflows can be constructed from operations of four different kinds.

Four separate execution layers, each distinguished by input/output types.





# 1 Deconstructing scientific workflows.

Modularity requires rules for how modules interact.

Theory: workflows can be constructed from operations of four different kinds.

Four separate execution layers, each distinguished by input/output types.



### 💊 🎇 👘 Jakob (Nordin alt. van Santen)

### Sample use case: Kilonova counterparts to gravitational wave events



The AMPEL alert archive, integrated astronomical catalogs and parallelizable T2 computations allow large data volume to be scanned.



JvS





Implements operations as units bound to a Tier.

Collects a series of units into a channel.





#### Jakob (Nordin alt. van Santen)



Implements operations as units bound to a Tier.

Collects a series of units into a channel.

| Shaper A    |             | Filter B    |            | ТЛ |
|-------------|-------------|-------------|------------|----|
|             | Filter A    |             | Shaper B   |    |
|             | Combiner I  |             |            | T1 |
|             |             | Combiner II |            |    |
| Augmenter I |             | Analyser I  |            | Т2 |
|             | Augmenter 2 |             | Analyser 2 | 12 |
| (           |             | Reaction 2  |            | TO |
| Reaction 1  |             |             |            | 13 |
|             |             |             |            |    |



#### Jakob (Nordin alt. van Santen)



Unit: python module inheriting from Ampel-interface class.

#### 

#### Class methods determine I/O.



### Jakob (Nordin alt. van Santen)

AMPEL Core:

- Workers at each tier executes units with requested input, allowing system control and parallelisation.
- Results stored in NoSQL (Mongo) DB.
- Built-in provenance tracking (event journal, logs and jobs)

Execute a job:

- locally to develop,
- at a cluster for large-volume archive runs
- in a live instance to analyze real-time data

\*Contrib User designed units Ampelcore

A.-interface

A.-alerts

A.-photometry

A.-ZTF

Jakob (Nordin alt. van Santen)

### Sample use case: Searching for gravitationally lensed supernovae



Exceedingly rare cosmological probes.



Extension of an existing unit allowed current pipelines to be reused for a lens searches.



### 3 States: key to provenance and deduplication

A state is the collection of data points associated with an (assumed) object at a specific time and visible to a specific user.

AMPEL systematically records every transient state as an immutable object.

Every computation is directly tied to a state, a key component of the enforced data provenance model. Also enables deduplication - identical computations on the same state are elided.

| (1) ObjectId("63790bc419a1a4c4e62cecd | 2") { 7 fields }                     |
|---------------------------------------|--------------------------------------|
|                                       | ObjectId("63790bc419a1a4c4e62cecd2") |
| 💻 link                                | 1395918359                           |
| stock                                 | 44361271                             |
| Channel                               | [ 2 elements ]                       |
| III dps                               | [ 44 elements ]                      |
| 🕨 💷 meta                              | [ 1 element ]                        |
| 🕨 💷 tag                               | [ 1 element ]                        |

|           | 3 (4)                             | ObjectId("63790bc419a1a4c4e62cea23")   | { 10 fields }  |
|-----------|-----------------------------------|--|--|
|           |                                   | _id  | ObjectId("63790bc419a1a4c4e62cea23")   |
|           | #                                 | config   | -741625996551693404  |
|           | #                                 | link   | 1395918359   |
|           | #                                 | stock  | 44361271   |
|           |                                   | unit   | T2XgbClassifier  |
| •         |                                   | channel  | [ 2 elements ]   |
|           | #                                 | code   | 0  |
| •         |                                   | meta   | [ 2 elements ]   |
| •         | 1                                 | tag  | [ 1 element ]  |
| -         |                                   | body   | [ 1 element ]  |
|           |                                   |  |  |
|           |                                   |  |  |
| • 🖸       | 8)                                | ObjectId("63790bc419a1a4c4e62cea27")   | { 10 fields }  |
| • ()      | 8)                                | Objectld("63790bc419a1a4c4e62cea27")<br>_id  | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")  |
| •         | 8)                                | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config  | { 10 fields }<br>Objectld("63790bc419a1a4c4e62cea27")<br>4519816789791135494   |
| ▼ €2      | 8)<br>(8)<br>#                    | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link  | { 10 fields }<br>Objectid("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359   |
| •         | 8)<br>(8)<br>#                    | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock   | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271   |
| ▼ 🖸       | ) (8)<br>#<br>#<br>#              | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock<br>unit                                   | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271<br>T2ElasticcReport   |
| •         | 8)<br>#<br>#<br>#<br>**           | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock<br>unit<br>channel                        | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271<br>T2ElasticcReport<br>[ 2 elements ]   |
| •         | 2 (8)<br>#<br>#<br>#<br>#<br>**   | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock<br>unit<br>channel<br>code                | { 10 fields }<br>Objectid("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271<br>T2ElasticcReport<br>[ 2 elements ]<br>0                                    |
| ▼ €⊇<br>} | (8)<br>#<br>#<br>#<br>**          | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock<br>unit<br>channel<br>code<br>meta        | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271<br>T2ElasticcReport<br>[ 2 elements ]<br>0<br>[ 5 elements ]                  |
| ▼ €2      | (8)<br>#<br>#<br>#<br>#<br>*<br>* | ObjectId("63790bc419a1a4c4e62cea27")<br>_id<br>config<br>link<br>stock<br>unit<br>channel<br>code<br>meta<br>tag | { 10 fields }<br>ObjectId("63790bc419a1a4c4e62cea27")<br>4519816789791135494<br>1395918359<br>44361271<br>T2ElasticcReport<br>[ 2 elements ]<br>0<br>[ 5 elements ]<br>[ 1 element ] |

### Sample use case: Photometric classification for LSST

LSST has immense science potential - *if* transients can be photometrically classified.



Modularity allowed AMPEL to be quickly adopted to the ELASTICC data challenge and emulate a realistic transient program.

Each classification tied to a state.



### 💊 🎇 👘 Jakob (Nordin alt. van Santen)

### Useful tools

Contains product endorsements

# Type annotations + mypy

- (Optionally!) declare types for variables, arguments, return values in Python >= 3.5
- Check for type safety (and stupid typos) with <u>mypy</u>, e.g. in refactoring large codebases with incomplete test coverage
- Bonus: more powerful tab completion in e.g. VSCode (Pylance) or vim (<u>Jedi</u>)

```
def get item(target: dict[str, int], key: str) -> None | int:
    return target.get(key)
```

```
def do_a_thing_with_a_list(target: list) -> None:
    return target.keys()
```

```
if name == " main ":
   get_item({"foo": 1}, "foo")
   get_item([("foo", 1)], "foo")
   get_item({"foo": 1.}, "foo")
```

Code with type annotations (and errors!)

```
> mypy mypy_demo.py
mypy_demo.py:5: error: "List[Any]" has no attribute "keys"
[attr-defined]
mypy_demo.py:9: error: Argument 1 to "get_item" has
incompatible type "List[Tuple[str, int]]"; expected "Dict[str,
int]" [arg-type]
mypy_demo.py:10: error: Dict entry 0 has incompatible type
"str": "float"; expected "str": "int" [dict-item]
Found 3 errors in 1 file (checked 1 source file)
```

#### Errors found statically by mypy



#### Jakob (Nordin alt. van Santen)

# pydantic

- Sometimes you just want a struct that
  - a. Generates \_\_init\_\_ and \_\_repr\_\_ for you => dataclasses
  - b. Can also load itself from an untrusted source, and give useful error messages when it fails => <u>pydantic</u>
- Annotations define expected types
- Basic validation with constrained types, e.g. "list with at least one item", "float < 0"</li>
- Arbitrary validation with <u>@validator</u>
- Validation logic stays with the data model

#### A basic pydantic model

```
from typing import Sequence
from pydantic import BaseModel
class ChannelModel(BaseModel):
    channel: int | str
    version: None | int | float | str
    active: bool = True
    members: None | Sequence[str]
if __name__ == "__main__":
    ChannelModel()
```

#### A validation error

| pydantic.error_wrappers.ValidationErr |
|---------------------------------------|
| or: 1 validation error for            |
| ChannelModel                          |
| channel                               |
| field required                        |
| (type=value_error.missing)            |
|                                       |

### Jakob (Nordin alt. van Santen)

### FastAPI

Does what it says: write REST APIs, fast

- Most of a REST API is input validation •
- **FastAPI** delegates validation to pydantic, and returns an informative 422 Unprocessable Entity on failure
- Declare handlers with decorators

```
async def get process(process: str) -> ProcessModel:
  for tier in range(4):
           doc = context.config.get(f"process.t{tier}.{process}", dict,
raise exc=True)
      raise HTTPException(status code=404, detail=f"{process} not found")
```

A GET handler with a single parameter

| •••  | $\langle \rangle = 0$  | ampel.zeuth                                 | ien.desy.de               | ମ୍ବର ଓ    | ₾       |
|--|--|---|---------------------------|-----------|---------|
|  | ert Arcl   | hive Se                                     | rvice                     | 3.1.0     | OAS3    |
| /api/ztf/archive/v3/ope<br>Query ZTF alerts is                                 | ssued by IPAC  |   |                           |           |         |
| Authorizat   | tion   |   |                           |           |         |
| Some endpoints re<br>token using the "A  | quire an authorizati<br>rchive tokens" tab o   | ion token. You can c<br>in the Ampel dashbo | reate a ZTF<br>ard. These | archive a | ccess   |
| and associated wit   | h your GitHub userr  | name.                                       |                           |           | poroioi |
| and associated wit Servers /api/ztf/archive/v                                  | h your GitHub userr  | name.                                       |                           | Authorize |         |
| and associated wit<br>Servers<br>/api/ztt/archive//<br>alerts Retri            | h your GitHub userr<br>v3 ~<br>eve alerts  | name.                                       |                           | Authorize |         |
| and associated with<br>Servers<br>/api/ztt/archive/<br>alerts Retria           | h your GitHub userr<br>v3 v<br>eve alerts<br>exts Post Alert Chu                     | unk   |                           | Authorize |         |
| and associated wit<br>Servers<br>/api/ztt/archive/<br>alerts Retri<br>POST /az | h your GitHub userr<br>v3 ~<br>eve alerts<br>erts Post Alert Chu<br>ert/(candid) Get | unk<br>Alert                                | ,                         | Authorize |         |

Autogenerated interactive API docs

### Zeuthen Data Science Seminar



#### Jakob (Nordin alt. van Santen)

# Authentication and Authorization

- Ampel is a multi-user system with a (read-only) API => restrict views to authorized users
- Most people have a GitHub account
   => use GitHub as identity provider
- <u>Authorization backend</u> issues a signed token containing username and org/group memberships
- Channel definitions include identities that should have access to data and logs

Example channel config (excerpt)



Decoded JWT payload in jwt.io debugger



## Renovate

- <u>Renovate</u> is a dependency updater with support for nearly every manager you can think of (think Dependabot, but more configurable)
- Example from our config:
  - a. Ampel sub-projects: update ASAP
  - b. Dev dependencies: bundle and update once a month, automerge if CI passes
  - c. Major updates: open a PR for each
- Run hosted or on-prem (e.g. for GitLab)

| Merged renovate  | eps): update depe<br>merged 1 commit int                           | endency m<br>master fr                 | om renovate                  | #210<br>/dev-depend       | lenc            |
|--|--|--|------------------------------|---------------------------|-----------------|
| renovate bot o   | commented on Ma  | r 13                                   |                              | Cor                       | tributor ···    |
| անվիք Ֆիսիստ   | 🕐 MEN  | DRenov                                 | ate                          |                           |                 |
| his PR contains the f  | ollowing updates:  |  |                              |                           |                 |
| Package  | Change   | Age                                    | Adopti<br>on                 | Passi<br>ng               | Confide<br>nce  |
| mypy (source,<br>changelog)  | ^1.0 -><br>^1.1  | 😂 41d                                  | 7%                           | <b>ð</b> 91%              | 🐥 high          |
|  |  |  |                              |                           |                 |
| onfiguration<br>Schedule: Branch<br>Itomerge - At any tir<br>Automerge: Enabl<br>Rebasing: Whene<br>teckbox. | creation - "every<br>ne (no schedule d<br>led.<br>ver PR becomes c | 1 month or<br>efined).<br>onflicted, o | n Monday be<br>or you tick t | efore 5am"<br>he rebase/i | (UTC),<br>retry |



## Miscellaneous morsels

- Stop arguing about style and just use <u>black</u>.
- Write tests without the boilerplate using <u>pytest</u>.
- Dependency lock files are a great way to keep a known-working set of dependencies for testing. Use a manager like <u>poetry</u>, <u>pdm</u>, or <u>conda-lock</u> to create and maintain them.
- Solving and installing large environments with <u>conda</u> is surprisingly slow. <u>micromamba</u> is much faster, and comes as a statically-linked binary.



### Lessons learned

# NoSQL databases are great!

- ...until you have your first many-to-many relationship
- SQL: store single copy, reference by id
- NoSQL: either
  - a. Keep multiple copies, forget about consistency
  - b. Build your own relational logic at the application level, using only idempotent and commutative operations
- Still a decent choice when you never update or delete records, or just have piles of unrelated documents.
- If you just want to avoid defining a schema, consider jsonb columns in PostgreSQL.



# Tightly-coupled Git repositories

- Sometimes you have lots of tightly-coupled code that you want to update atomically, but want different permissions for each component (some astronomers care about their code being secret)
- Git makes this so much harder than SVN. Can enforce write permissions with PRs and CI, but there is no way to make fine-grained read permissions
- Least bad workaround so far:
  - a. split into multiple repos
  - b. publish packages from CI for basically any change
  - c. update dependency spec on downstream packages
- But: can't detect when a change breaks a downstream package

### Jakob (Nordin alt. van Santen)

# Enabling scientific software development is not enough

- AMPEL was designed to encourage and support continuous scientific software development.
- ... this does not mean it happens.

Multiple cases of well working units constructed by knowledgeable scientists. Appreciated by science community, but not supported once v1 is running.



### Labels matter

- It is hard to introduce new workflow concepts.
- Broker is a new thing everything is a broker what is a broker?
- Code-To-Data actually seems to work in communication:

AMPEL allows scientific time-domain software to be locally developed, shared with other users and uploaded to high performance computer centers - Code-To-Data.

