

GRAPHS REDUX

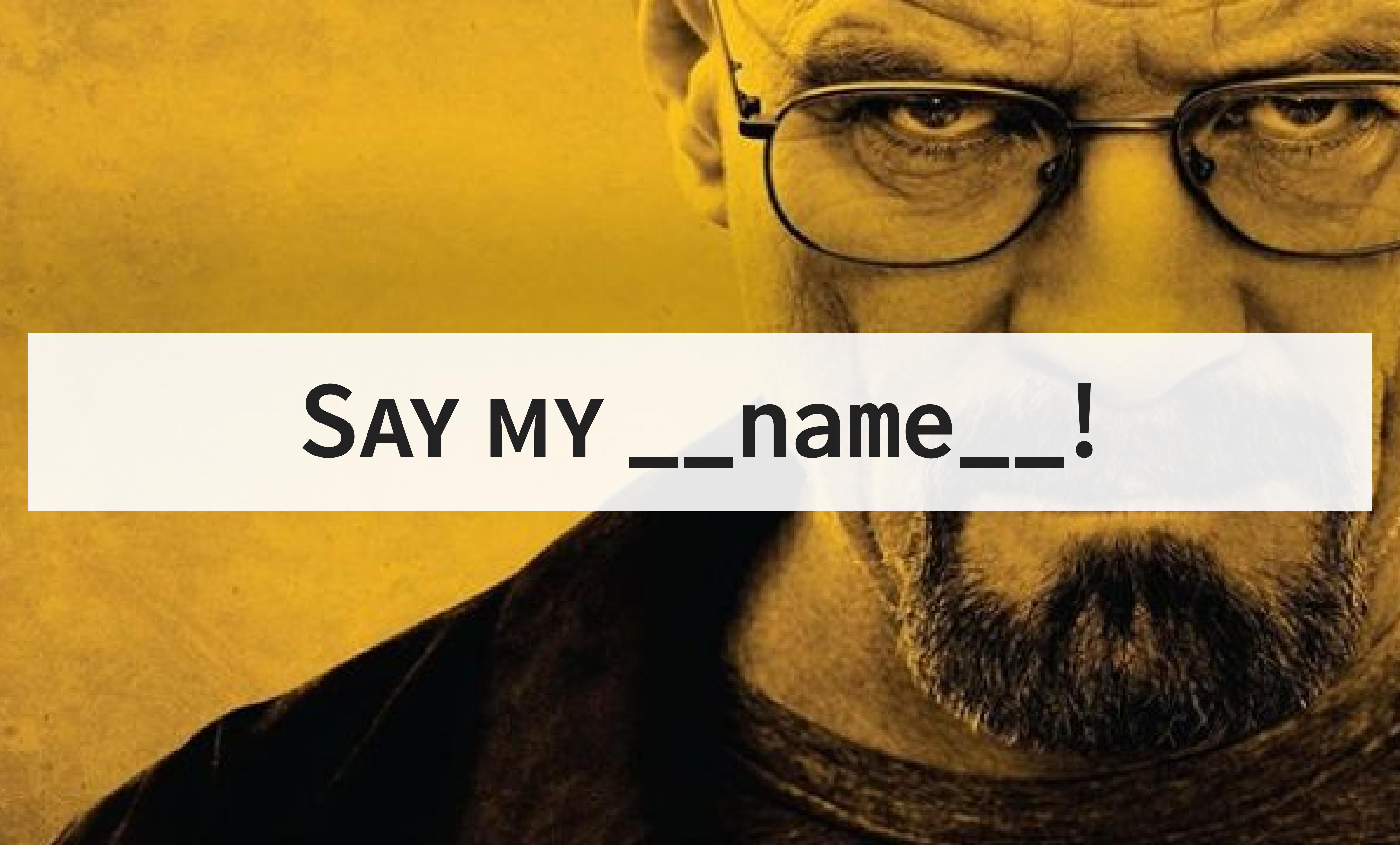
Dr. Scott Gorlin

Harvard University

Fall 2018

AGENDA

- Say my __name__!
- Data Pipeline
- Advanced Luigi
- Where We Are
- Dask
- Parquet
- Readings



SAY MY __name__!

A MISUNDERSTOOD IDIOM

```
if __name__ == '__main__':  
    ...
```

This may not be what you want!

WHAT'S IN A `__name__`?

Files:

```
main.py  
package/  
    __init__.py  
module.py
```

`python main.py:`

```
if __name__ == '__main__':  
    import sys  
    print(sys.path[0]) # /Users/scott/repos/package  
    print(__name__) # __main__  
    import package  
    from package import module  
    print(package.__name__) # package  
    print(module.__name__) # package.module
```

Looks ok so far!

WHAT'S IN A `__name__`?

package.module:

```
import logging

logger =
logging.getLogger(__name__)

if __name__ == '__main__':
    import sys
    print(sys.path[0])
    import package
```

python package/module.py:

```
/Users/scott/repos/package/package # !!!
Traceback (most recent call last):
  File "package/module.py", line 7, in <module>
    import package
ModuleNotFoundError: No module named 'package'
```

Python changed the python path, and altered the working code!

module is treated like a top level module, not part of a package

WHAT'S IN A `__name__`?

package.module:

```
import logging

logger =
logging.getLogger(__name__)

if __name__ == '__main__':
    import sys
    print(sys.path[0])
    import package
```

Other programmatic usages of `__name__` will break!

eg anything any other from package.module

```
import logger
```

DO NOT DIRECTLY RUN PACKAGES

Never directly run `python some_package/__init__.py` or any python file that is part of a library. It is not designed for that.

Your IDE can trick you: never ‘debug package/module.py’ directly; use a scratch file, `main.py`, `junk.py`, etc at the top of the repo

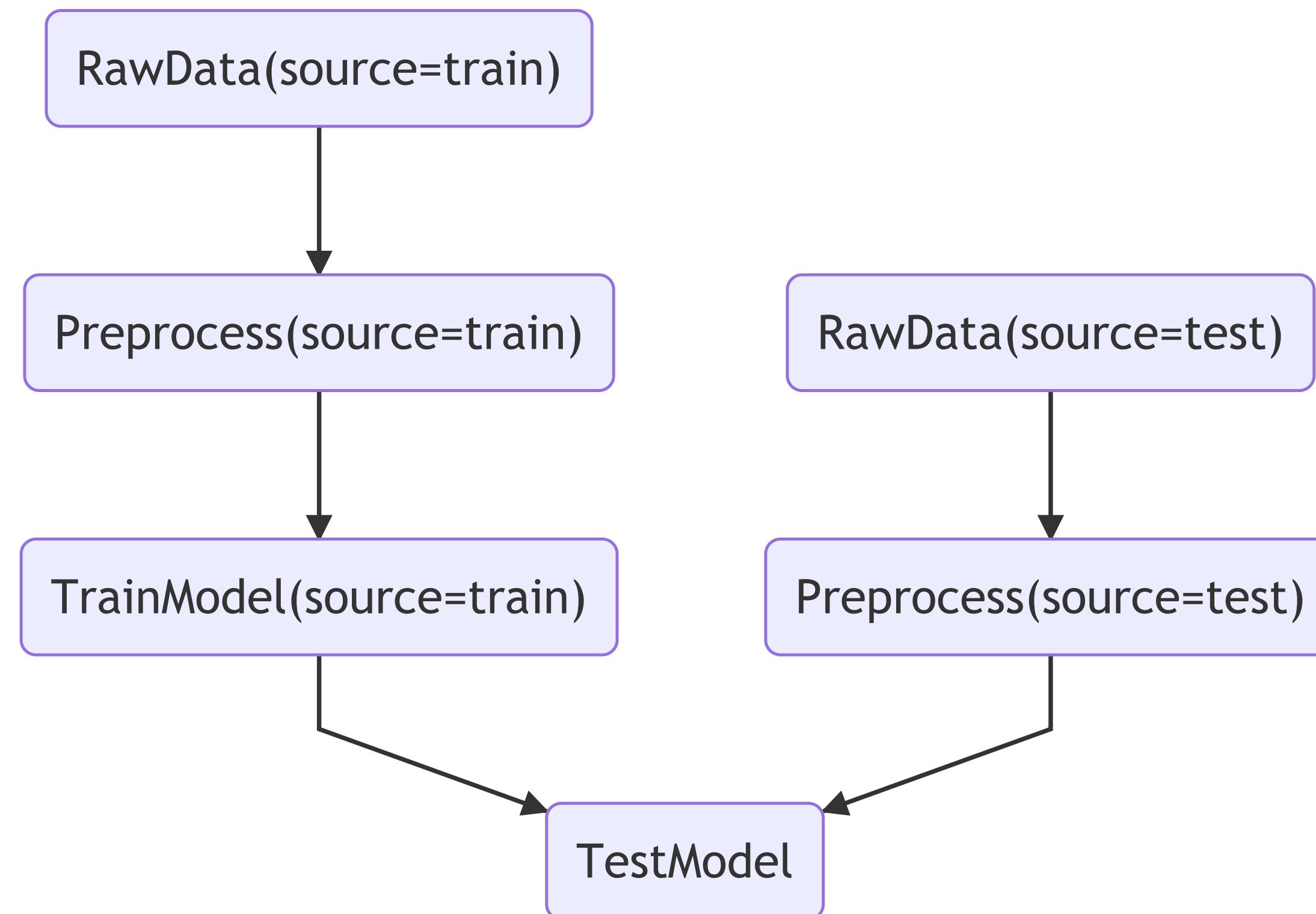
You can directly run a package with `python -m package` if you implement `package/__main__.py` in a similar form.



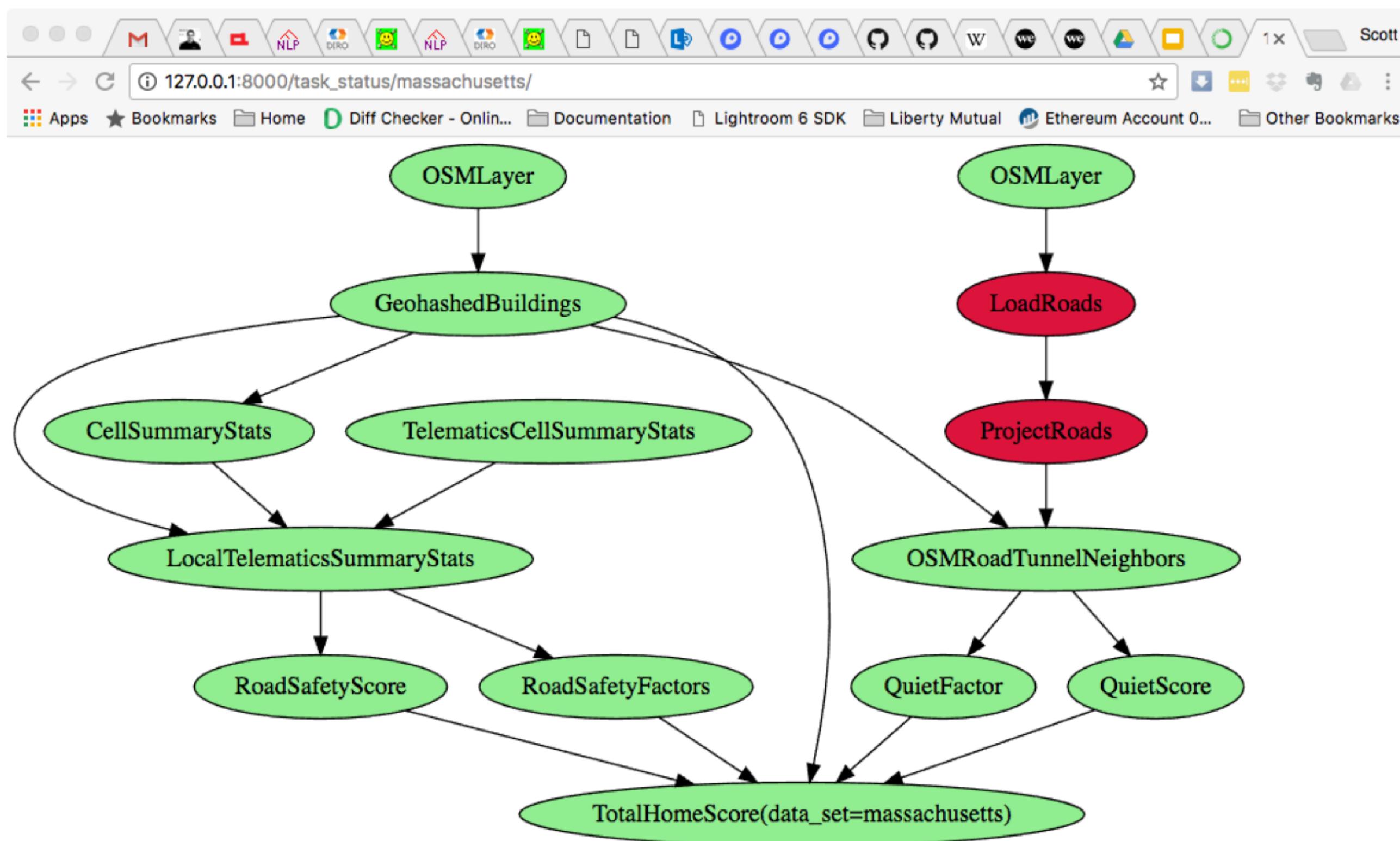
DATA PIPELINE

THE BIG PICTURE

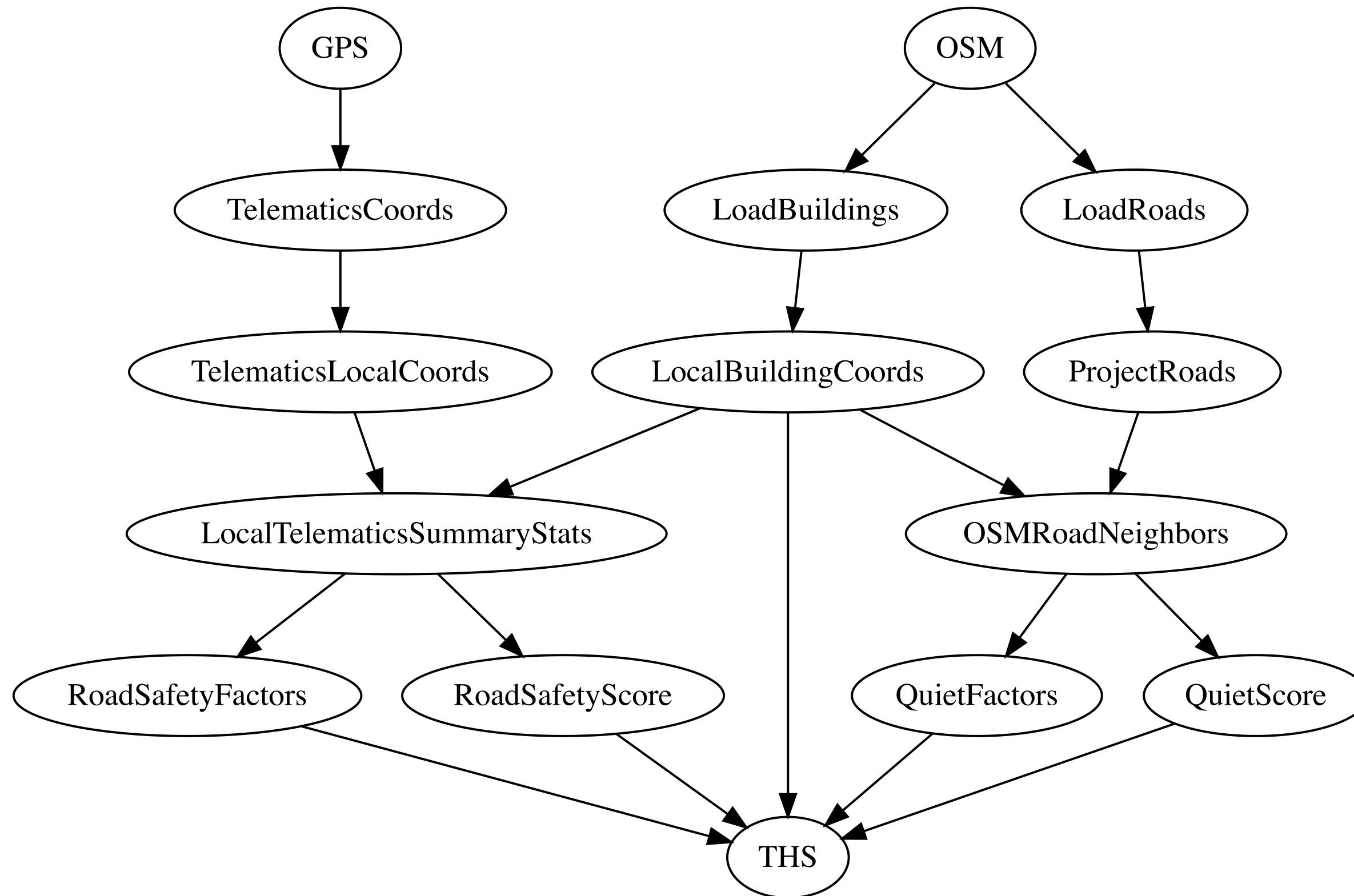
Using Task's, we now can see our entire project from above



THE BIG PICTURE

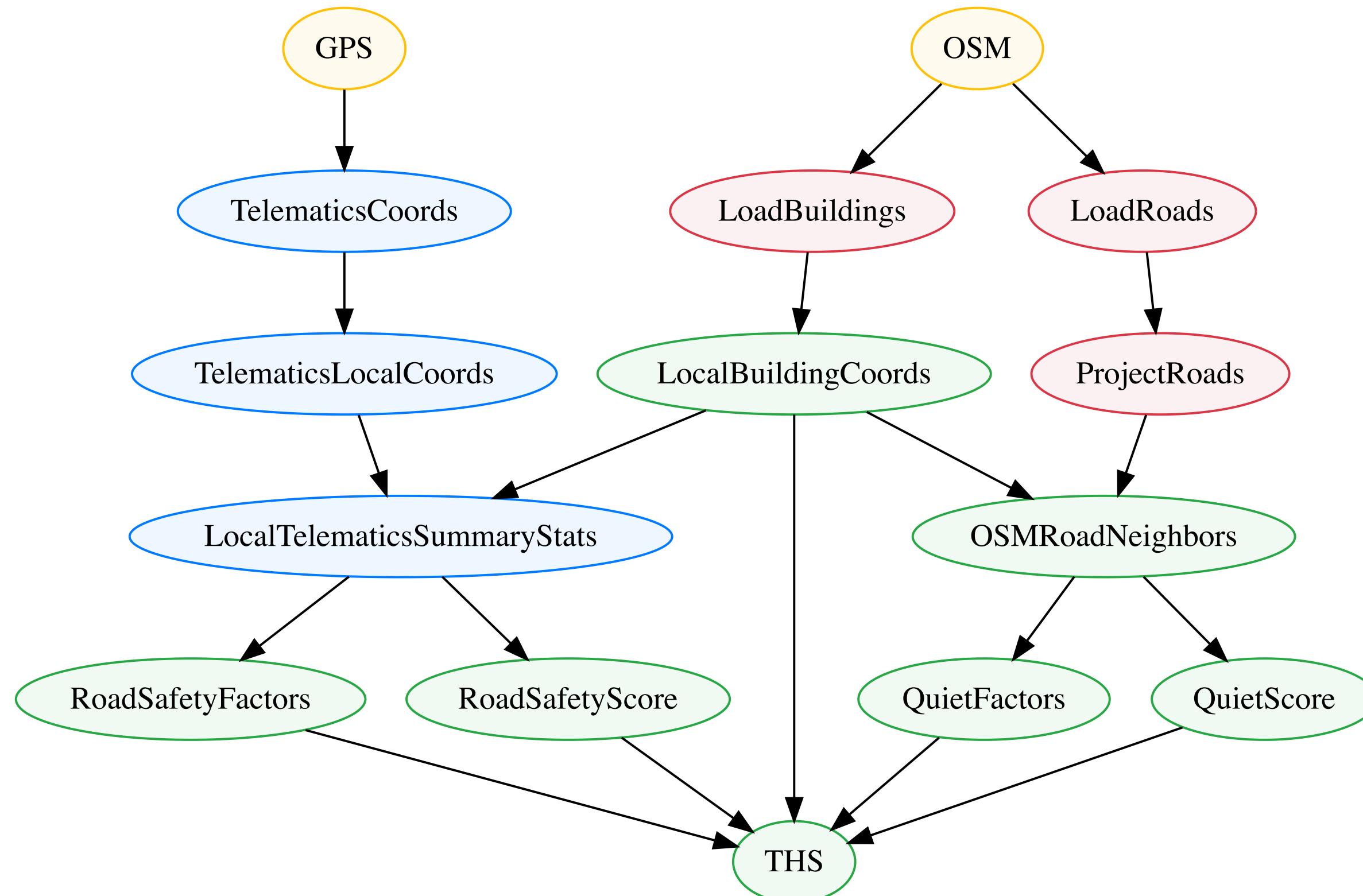


THE BIG PICTURE



Total Home Score

THE BIG PICTURE



Spark/Parquet

Dask/Parquet

External

Django/SQL

SCHEDULING

LUIGI

REQUIREMENTS

Including dynamic dependencies

OUTPUTS

Including logic, salting, etc

EXECUTORS

Running everything

DELEGATED: .run()

Just pass the target paths!

SPARK

Kick off MR/ETL

OOZIE

Another scheduling system!

...

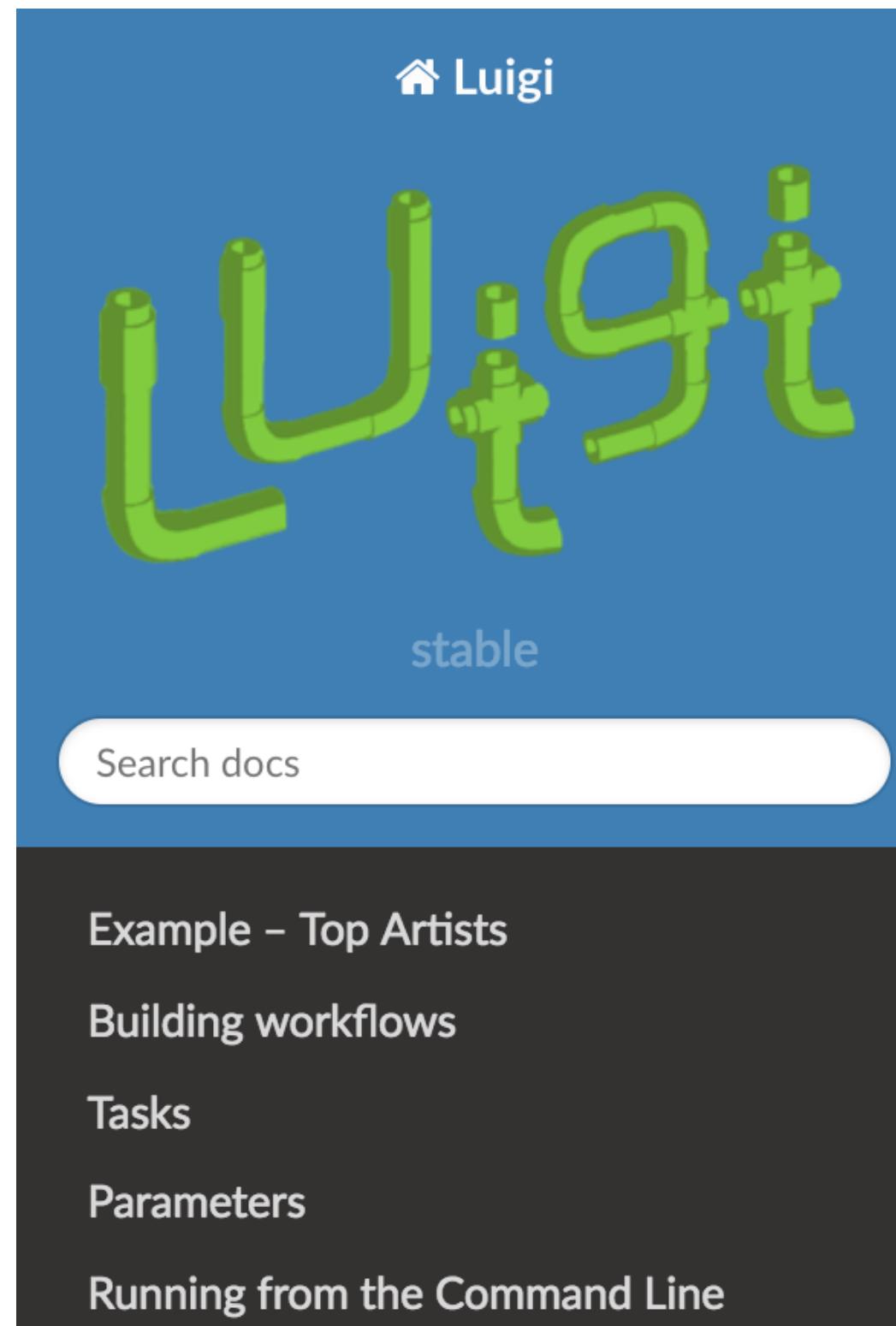


ADVANCED LUIGI

TASK TYPES

So far we've seen Task and ExternalTask, but many more exist!

Luigi has a rich set of contributions, mostly for running non-python tasks



Submodules

- [luigi.contrib.batch module](#)
- [luigi.contrib.bigquery module](#)
- [luigi.contrib.bigquery_avro module](#)
- [luigi.contrib.dataproc module](#)
- [luigi.contrib.docker_runner module](#)
- [luigi.contrib.ecs module](#)
- [luigi.contrib.esindex module](#)
- [luigi.contrib.external_program module](#)
- [luigi.contrib.ftp module](#)
- [luigi.contrib.gcp module](#)
- [luigi.contrib.gcs module](#)
- [luigi.contrib.hadoop module](#)
- [luigi.contrib.hadoop_jar module](#)
- [luigi.contrib.hive module](#)
- [luigi.contrib.kubernetes module](#)
- [luigi.contrib.lsf module](#)

EXTERNAL PROGRAMS

```
# luigi.contrib.external_program
class ExternalProgramTask(Task):
    def program_args(self):
        raise NotImplementedError()
    def run(self):
        args = self.program_args()
        logger.info(
            'Running command: %s',
            '\n'.join(args))
        subprocess.check_call(args)
```

```
# luigi.contrib.spark
class SparkSubmitTask(ExternalProgramTask):
    def program_args(self):
        return [
            'spark-submit',
            '--name', self.name,
            self.app
            ...
        ]
class PySparkSubmitTask(SparkSubmitTask):
    ...
```

EXTERNAL PROGRAMS

Your luigi target doesn't control how an external program writes data - it can only point to the data once written.

You should choose a target that matches the atomicity pattern of the program.

Eg, luigi.contrib.S3FlagTarget which looks for a _SUCCESS flag left by many hadoop/EMR jobs, including spark.

EXTERNAL PROGRAMS

Programs come in two flavors:

LOCAL/WRAPPED RUNNER

The job process starts and locally succeeds or fails

Examples: `python main.py`,
`spark --deploy-mode client`

DISPATCHED RUNNER

Process submits a remote job and successfully terminates

The task output could be job id; next task queries job status
(out of scope for this class)

Examples: Oozie, EMR step,
`spark --deploy-mode cluster`

THE INEVITABLE CREEP

TASK EXECUTOR TYPES	OUTPUT TYPES	VARIANTS
<ul style="list-style-type: none">• External tasks• Spark tasks• SQL tasks• Kube tasks• Docker tasks	<ul style="list-style-type: none">• Local targets• S3/hdfs Targets• Partitioned Targets• DB Targets	<ul style="list-style-type: none">• Salted targets• Your task hierarchy

Your tasks have a high risk of multi inheritance!

TASK SKELETON

- Parameters
 - Already declarative
- Run
 - Seems like the main job of the task, keep as a method
- Requires, Output
 - Seem ripe for declarative composition

```
import luigi

class MyTask(luigi.Task):
    param = luigi.Parameter(default=42)

    def requires(self):
        return SomeOtherTask(self.param)

    def run(self):
        f = self.output().open('w')
        print >>f, "hello, world"
        f.close()

    def output(self):
        return luigi.LocalTarget('/tmp/foo/bar-%s.txt' % self.param)

if __name__ == '__main__':
    luigi.run()
```

The business logic of the task

Where it writes output

What other tasks it depends on

Parameters for this task

REQUIREMENTS

Requirements can be a single task, a list, or a dict

They often share some parameters with the task, and should be formed
with `task.clone(requirement_class, **param_overrides)`

REQUIREMENTS

We can implement this as a descriptor and a registry!

```
class Requirement:  
    def __init__(self, task_class, **params):  
        ...  
  
    def __get__(self, task, cls):  
        return task.clone(  
            self.task_class,  
            **self.params)
```

```
class Requires:  
    def __get__(self, task, cls):  
        return lambda : self(task)  
  
    def __call__(self, task):  
        # Search task.__class__ for Requirements  
        # return instances
```

```
class MyTask(Task):  
    # Replace task.requires()  
    requires = Requires()  
    other = Requirement(OtherTask)  
  
    def run(self):  
        # Convenient access here...  
        with self.other.output().open('r') as f:  
            ...  
  
>>> MyTask().requires()  
{'other': OtherTask()}
```

This will be on your next pset!

OUTPUT

Output tends to follow a few rules you can choose from:

1. Is there a relationship between the file path and the task name?
 - eg, {DATA_ROOT}/{param}/{task_name}.txt, with a configurable root directory
2. Is the output salted?
3. Do you follow defaults, like encrypting things?
4. Is output all local, S3 targets, etc? Do you need special authentication?

The output logic probably isn't all that related to your task classes, and likely can be solved using composition

OUTPUT

```
class LocalTargetOutput:  
    def __init__(self, file_pattern='{task}.txt'):  
        ...  
    def __get__(self, task, cls):  
        return lambda: LocalTarget(...)
```

```
class MyTask:  
    # Compose output  
    output = LocalTargetOutput()  
  
>>> MyTask().output()  
LocalTarget()
```



WHERE WE ARE

```
1 def primes(int nb_primes):
2     cdef int n, i, len_p
3     cdef int p[1000]
4     if nb_primes > 1000:
5         nb_primes = 1000
6
7     len_p = 0 # The current number of elements in p.
8     n = 2
9     while len_p < nb_primes:
10         # Is n prime?
11         for i in p[:len_p]:
12             if n % i == 0:
13                 break
14             else:
15                 p[len_p] = n
16                 len_p += 1
17
18         n += 1
19
20     # Let's return the result in a python list:
21     result_as_list = [prime for prime in p[:len_p]]
22
23     return result_as_list
```

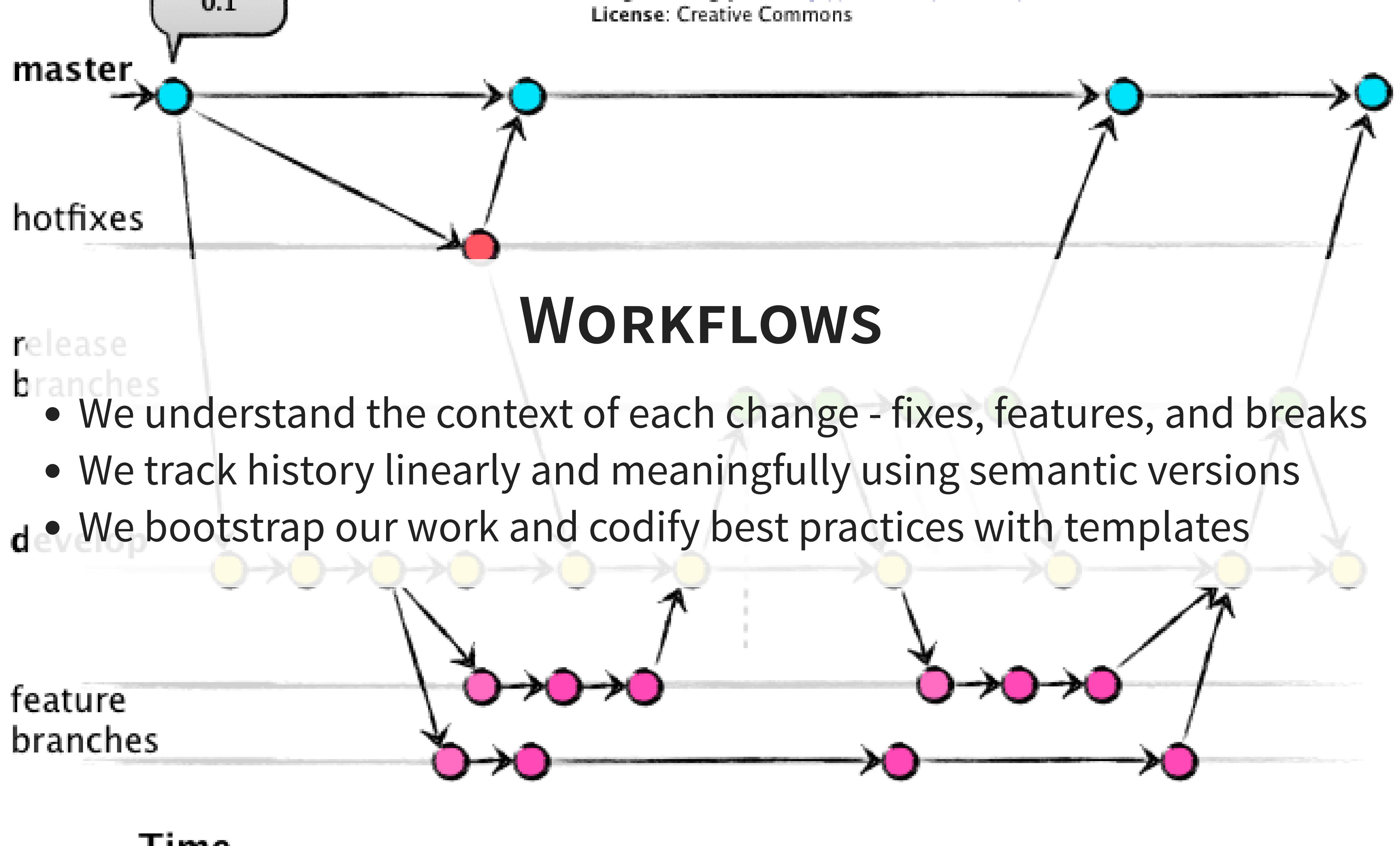
“PYTHON”

- We have repeatable and appropriately specified virtual environments
- We can choose between minimal reqs and frozen dependencies
- We understand the role of pipenv for library and app development



TESTING

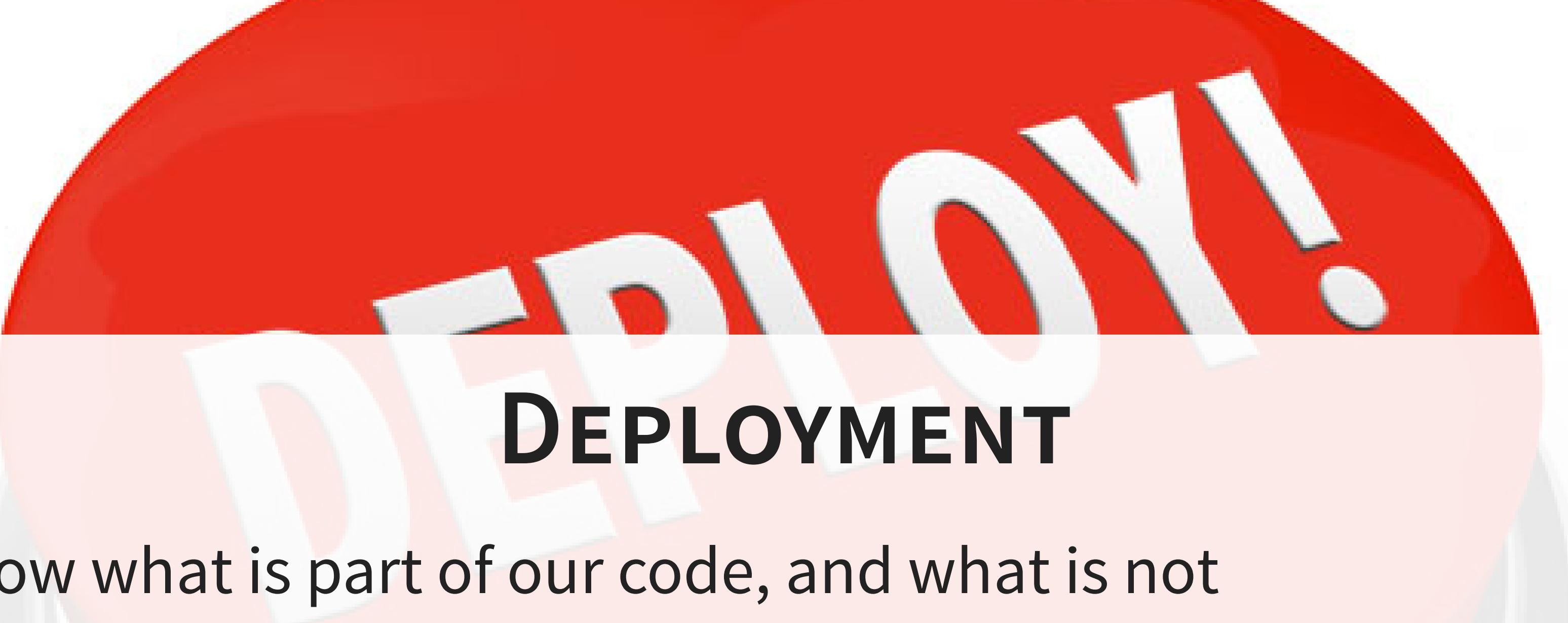
- We know the importance of unit testing, how to test, and what to test
- We can measure coverage, including code branches
- We know the environment matters, and test libraries where it counts
- We can mock out code, even faking integration with other systems





HIGHER LEVELS

- We recognize meta-patterns in code at levels higher than a function
- We wrap, register, and alter functions using decorators
- We provide reusable context to code using context managers



DEPLOYMENT

- We know what is part of our code, and what is not
- We configure our deployments with environment variables, ensuring our code is useable anywhere
- We handle data and secrets with discretion and privacy

LOOPING

- We see past for loops and recognize the higher looping primitive
- We can write stateless, functional code that expresses what we want, without telling the machine what to do
- We know the tradeoffs between efficiency, clarity, and diagnosability
- We know why we iterate, why we map, and why we reduce

FUNCTIONAL CODING

- We understand that our code is data, and may be operated on
- We know when state is valuable, and when it fails us
- We can encapsulate logic in functions that we pass as arguments to higher frameworks
- We strive to be declarative in all that we do

ASSEMBLY

PROCEDURAL

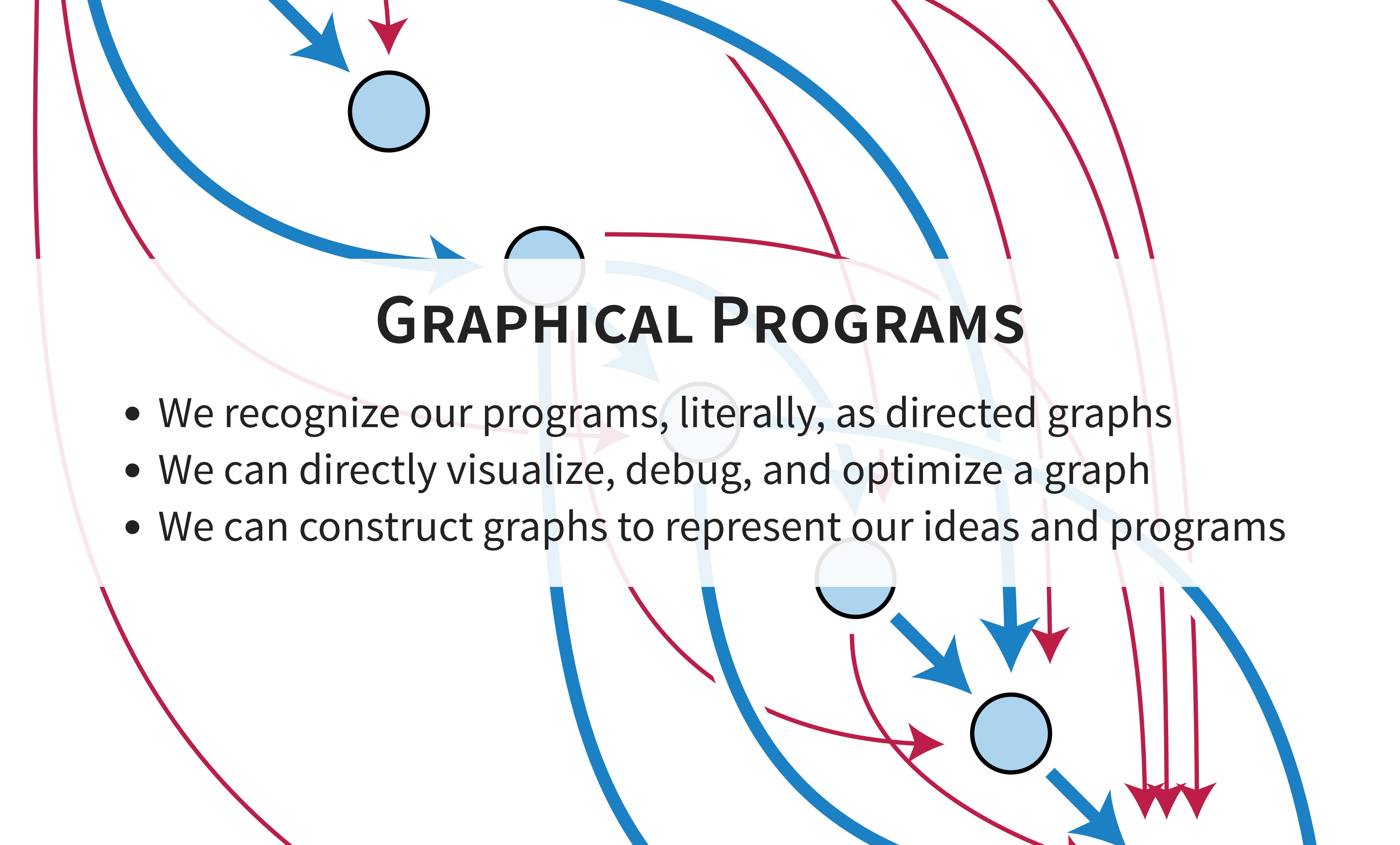
OBJECT ORIENTED

FUNCTI

COMPOSITION

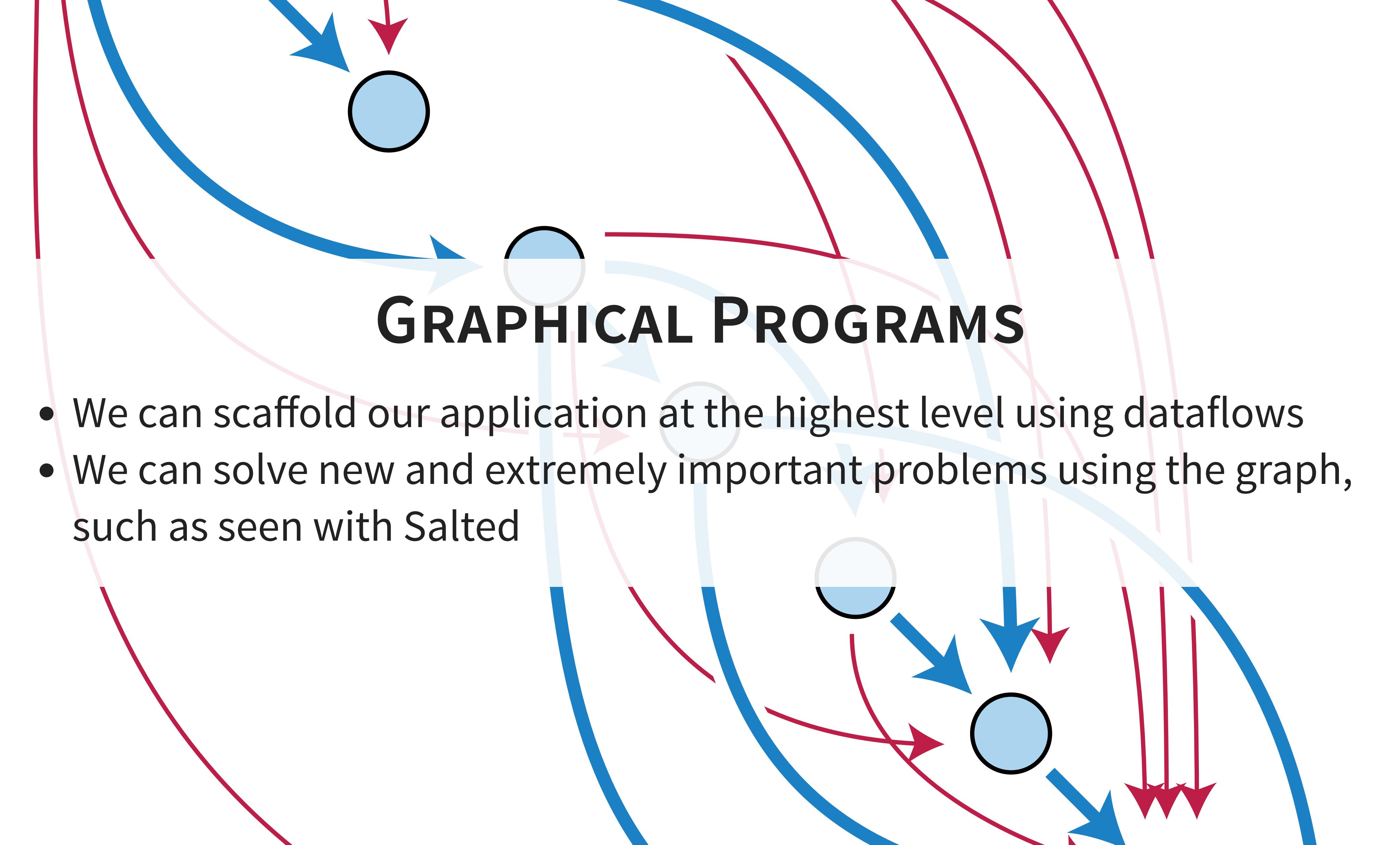
- We look for new ways to simplify and reuse our logic, such as mixins and composition
- We can add rich, reusable, and declarative properties with descriptors





GRAPHICAL PROGRAMS

- We recognize our programs, literally, as directed graphs
- We can directly visualize, debug, and optimize a graph
- We can construct graphs to represent our ideas and programs



GRAPHICAL PROGRAMS

- We can scaffold our application at the highest level using dataflows
- We can solve new and extremely important problems using the graph, such as seen with Salted



END OF LIST!

Midterm scope ends here

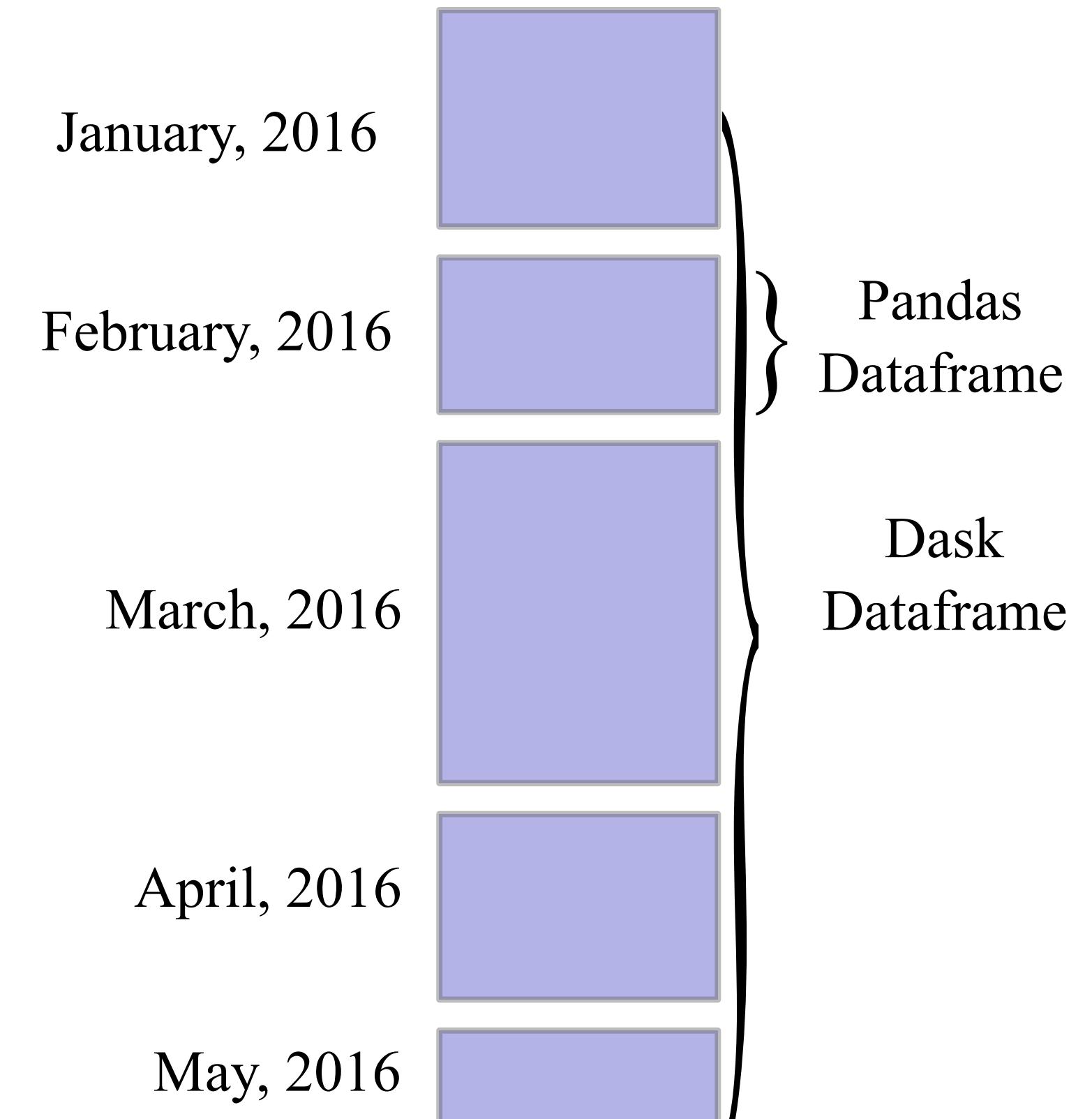
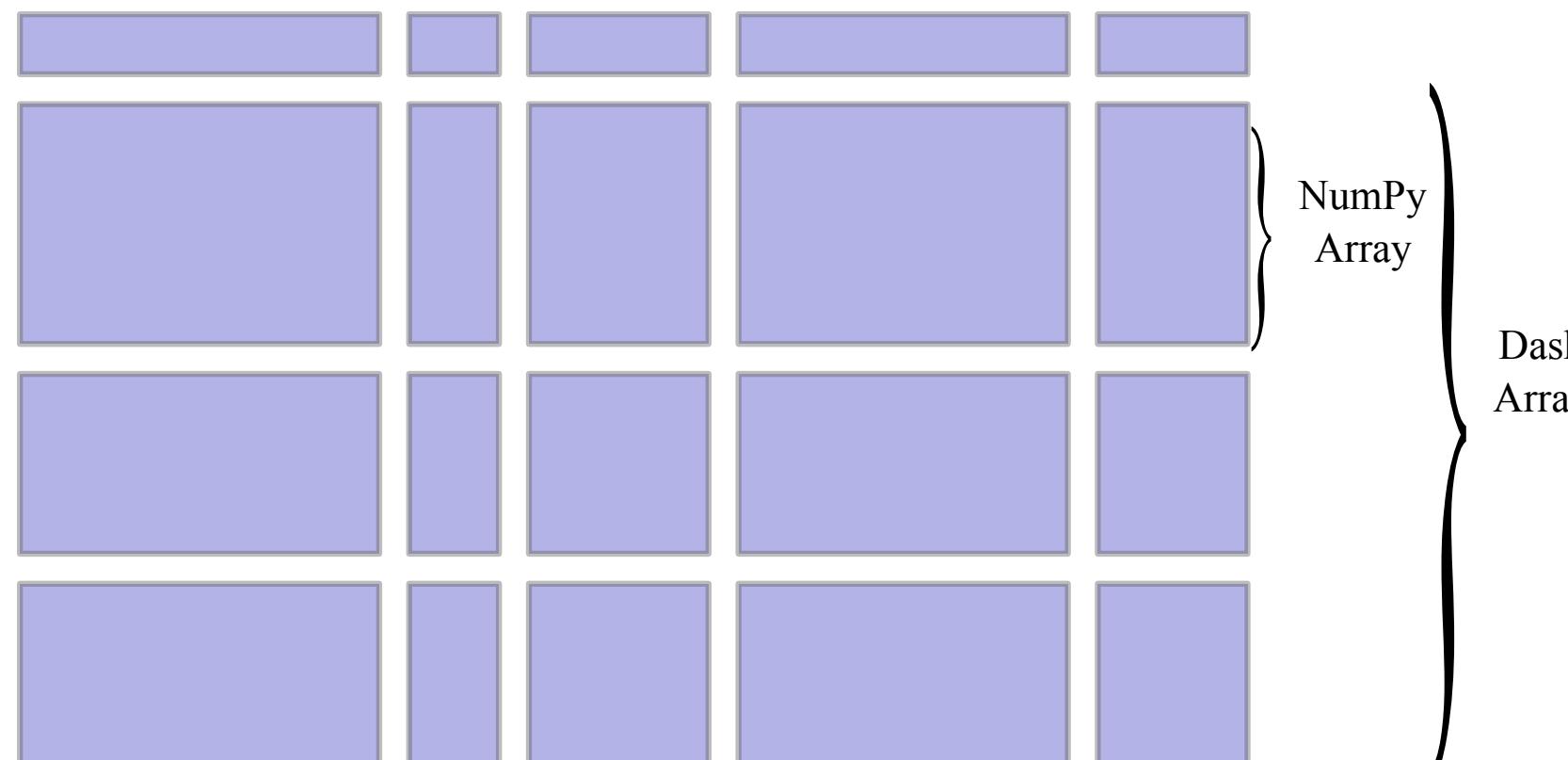


DASK

DASK

COMPUTATIONAL GRAPHS FOR PYTHON

- Out-Of-Memory numpy, pandas, text
- Chunked, distributed, delayed
- Can replace Spark for many workflows



dask.DataFrame

A partitioned, symbolic pandas frame

```
>>> import dask.dataframe as dd  
>>> df = dd.read_csv('2014-*csv')  
>>> df.head() # Implicit 'compute'  
   x  y  
0  1  a  
1  2  b  
2  3  c  
3  4  a  
4  5  b  
5  6  c
```

```
>>> df  
Dask DataFrame Structure:  
              x      y  
npartitions=3  
None          int64  object  
None          ...    ...  
None          ...    ...  
None          ...    ...  
Dask Name: from-delayed, 9 tasks
```

We don't notice it, but dask *has not* read the data yet. It did some brief checking to get row headers and column types, but data is not yet in memory.

dask.DataFrame

A partitioned, symbolic pandas frame

```
>>> pdf = pd.read_excel('hashed.xlsx')
>>> df = dd.from_pandas(pdf, chunkszie=10)
>>> df
Dask DataFrame Structure:
      learn project
npartitions=5
0d733444    object    object
33174d90        ...        ...
...            ...        ...
d9a774d1    object    object
ffc21800        ...        ...
Dask Name: from_pandas, 5 tasks
```

```
>>> df.loc['57019809']
Dask DataFrame Structure:
      learn project
npartitions=1
57019809    object    object
57019809        ...        ...
Dask Name: loc, 6 tasks
```

Notice the *predicate pushdown*.
This row index must fall within a
single partition, so dask only
returns that.

dask.DataFrame

Note the `_meta`!

```
>>> df = pd.DataFrame(  
...     {'a': [1, 2, 3],  
...      'b': ['x', 'y', 'z']})  
>>> ddf = dd.from_pandas(  
...     df, npartitions=2)  
>>> ddf._meta  
Empty DataFrame  
Columns: [a, b]  
Index: []  
>>> ddf._meta.dtypes  
a    int64  
b    object  
dtype: object
```

```
# Apply function to every partition  
ddf.map_partitions(  
    foo,  
    # Provide meta if necessary  
    meta=pd.DataFrame(  
        {'a': [1], 'b': [2]})  
    )  
) # cf ddf.apply()
```

Dask must know the shape of
your delayed data, and can't
always guess.

Provide it when you can

Row group 0

Column a

Page 0

Page header (ThriftCompactProtocol)

Repetition levels

Definition levels

values

Page 1

Column b

Row group 1

PARQUET

Footer

FileMetaData (ThriftCompactProtocol)

- Version (of the format)
- Schema
- extra key/value pairs

Row group 0 meta data:

Column a meta data:

- type / path / encodings / codec

number of values

offset of first data page

offset of first index page

- compressed/uncompressed size
- extra key/value pairs

column "b" meta data

Row group 1 meta data

Footer length (4 bytes)

COLUMN STORES

Row STORE

Vertical Partitions

id	name	grade
1	Scott	89
2	Andrew	91
99	Mary	87

COLUMN STORE

Vertical *and* horizontal partitions

id	name	id	grade
1	Scott	1	89
2	Andrew	2	91
99	Mary	99	87

Most operations only care about a few columns
Don't read them all if you don't need to!

COLUMN STORES

ANY TABLE...

A	B	C
A1	B1	C1
A2	B2	C2
A3	B3	C3

Row STORAGE

Any csv:

A1	B1	C1	A2	B2	C2	A3	B3	C3
----	----	----	----	----	----	----	----	----

COLUMN STORAGE

A better format:

A1	A2	A3	B1	B2	B3	C1	C2	C3
----	----	----	----	----	----	----	----	----

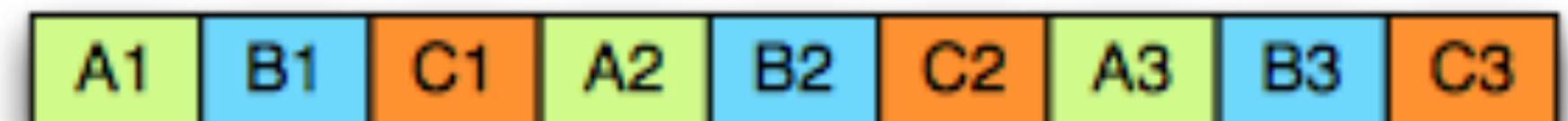
COLUMN STORES

Column stores offer superior:

- Compression, due to homogeneity
- Selective column reads
- IO, due to buffering

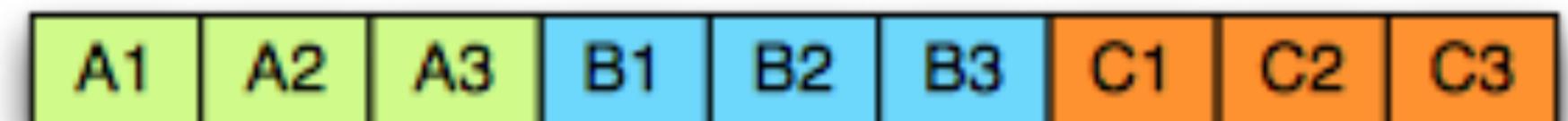
Row Storage

Any csv:



Column Storage

A better format:

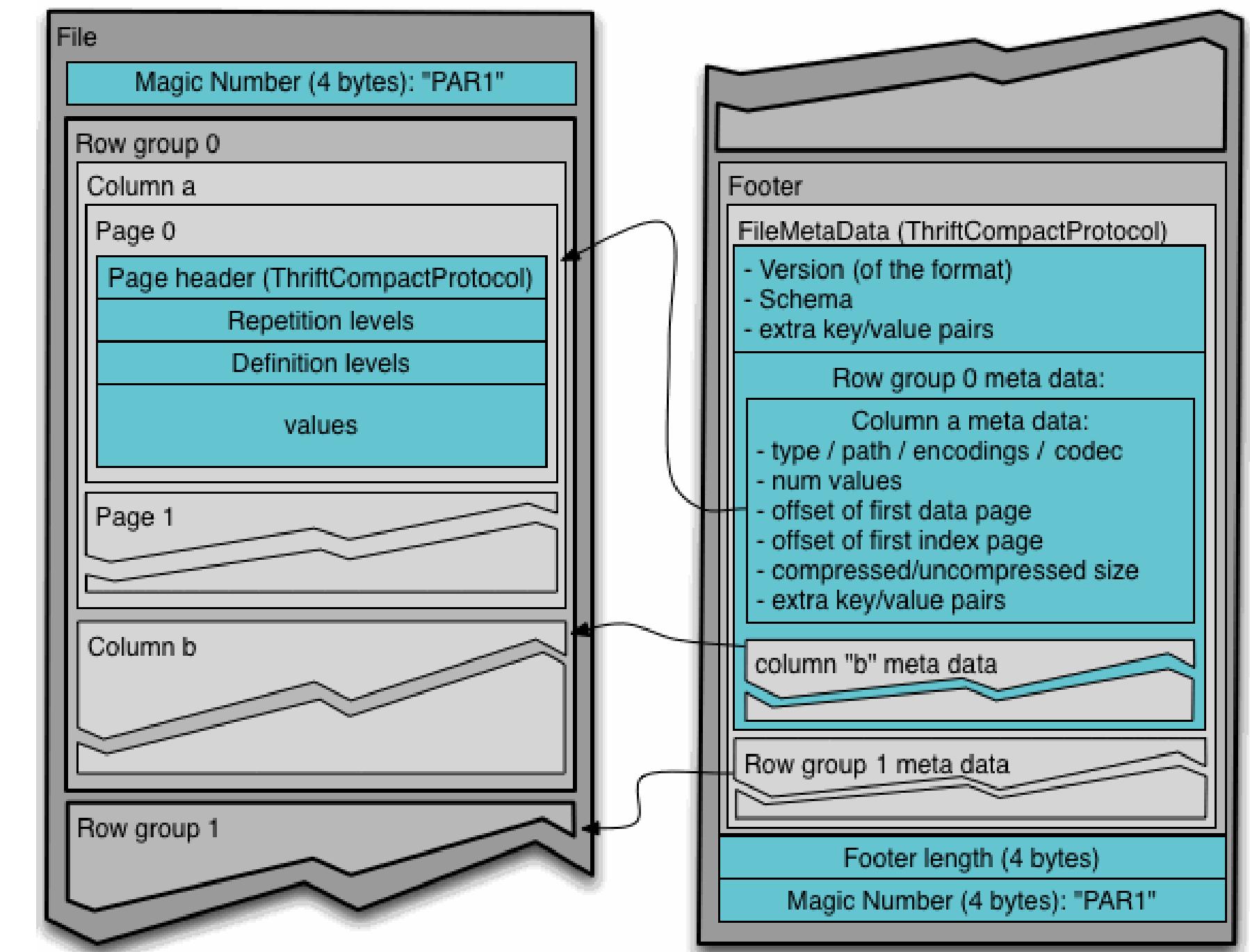


APACHE PARQUET

Emerging standard for distributed,
columnar binary data

Statistics for predicate pushdown

<http://parquet.apache.org/>



DATAFRAMES ♥ PARQUET

```
read_parquet(  
    path,  
    # Only read what you want!  
    columns=['a', 'b'],  
  
    # Filter what you read!  
    filters=[  
        ('a', '>', 5)  
    ],  
)
```

STATISTICS

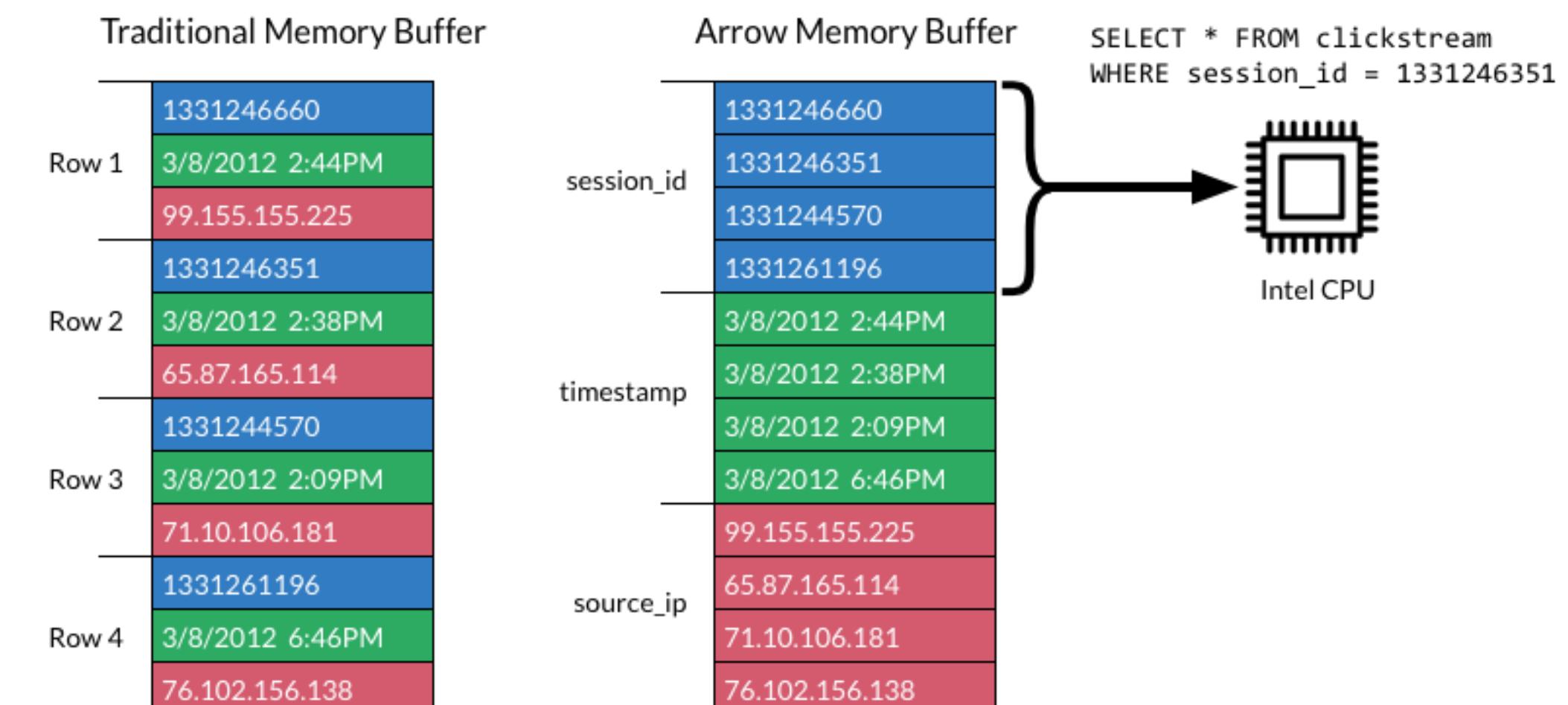
Parquet stores partition min/max values to allow for index and block filtering too!

APACHE ARROW

Unification of column stores - Pandas, Parquet, Spark, ...

<http://arrow.apache.org/>

	session_id	timestamp	source_ip
Row 1	1331246660	3/8/2012 2:44PM	99.155.155.225
Row 2	1331246351	3/8/2012 2:38PM	65.87.165.114
Row 3	1331244570	3/8/2012 2:09PM	71.10.106.181
Row 4	1331261196	3/8/2012 6:46PM	76.102.156.138



THE PLAYERS

PARQUET

Established on-disk
columnar standard
You will use directly
and extensively

ARROW

Emerging in-
memory standard
Basis for Pandas2
Integrated into
Parquet code base
Exciting, but you
likely won't touch
directly

FEATHER

POC minimal disk
format for arrow
We don't care

READINGS

For post midterm:

- Dask Overviews (esp DataFrame)
- Apache Arrow and the “10 Things I Hate About pandas”
- Apache Arrow and Apache Parquet
- Dremel made simple with parquet
- Skim slides: The Columnar Roadmap