

Up Your Bus Number

A Reproducible Data Science Workflow

Kjell & Amy Wooding



https://bit.ly/bus_number



Bus Number:

The number of people who need to get hit by a bus before your data science project becomes *irreproducible*.

This might be **zero**.

While you wait...

https://bit.ly/bus_number_help

hackalog / bus_number

Unwatch 1 Star 0 Fork 1

Code Issues 4 Pull requests 1 Projects 0 Wiki Insights Settings

Getting Started

Hackalog edited this page 2 hours ago · 2 revisions

Let's get you set up to do Reproducible Data Science.

The Bare Minimum

- + cookiecutter
- conda (and then python >= 3.6)
- make

Getting Started

1. Install [cookiecutter](#), then use it to install the `pydata_nyc` branch of `cookiecutter-easydata`:

```
cookiecutter https://github.com/hackalog/cookiecutter-easydata.git \
--checkout pydata_nyc
```

Pages 2

Find a Page...

Getting Started

Github Workflow Cheat Sheet

+ Add a custom sidebar

Clone this wiki locally

<https://github.com/hacka...>

https://bit.ly/bus_number

Who Needs Reproducible Data Science?

https://bit.ly/bus_number

Meet Bjørn

- Bjørn is a dot-com millionaire. Currently he heads the Ikea R&D kitchen in Sweden.
- Bjørn employs a large number of Finnish line cooks. He can't understand a word they say.
- Bjørn needs a **trained model** to do real-time translation from Finnish to Swedish.
- Even though test kitchens have notoriously high turnover, reproducibility means Bjørn can ask his sous-chef to **deploy updated models** whenever needed.



Source: <http://muppet.wikia.com/wiki/File:Swedish-chef.png>

Meet Mark

- Mark used to be an astronaut. After a rough year, he decided to change careers, and now works for a fashion magazine in NYC.
- Mark has to keep up with thousands of new fashions every week. More if it's Fashion Week.
- Mark wants to **publish a paper and software library** for accelerating fashion review using science.
- Reproducibility means his **peer reviewers** can reproduce his work, his paper will get accepted, his library will get pulled into scikit-learn, and he will transform the fashion industry forever.



Source: Youtube / 20th Century Fox

https://bit.ly/bus_number

Meet Annie

- Annie writes high frequency trading software. She moved to NYC after getting tired of the fast pace of life on the West Coast.
- Annie is constantly doing **EDA** on financial data in order to get that edge she can exploit for her firm.
- With thousands and thousands of analyses to sort through, there's no way Annie can remember everything she has done to her data.
- Reproducible data science gives Annie a way to **keep track** of the data she has worked with, the models she has built, and the experiments she has run, so that successful experiments can make it to the trading floor.



Source: <https://www.imdb.com/title/tt0111257>

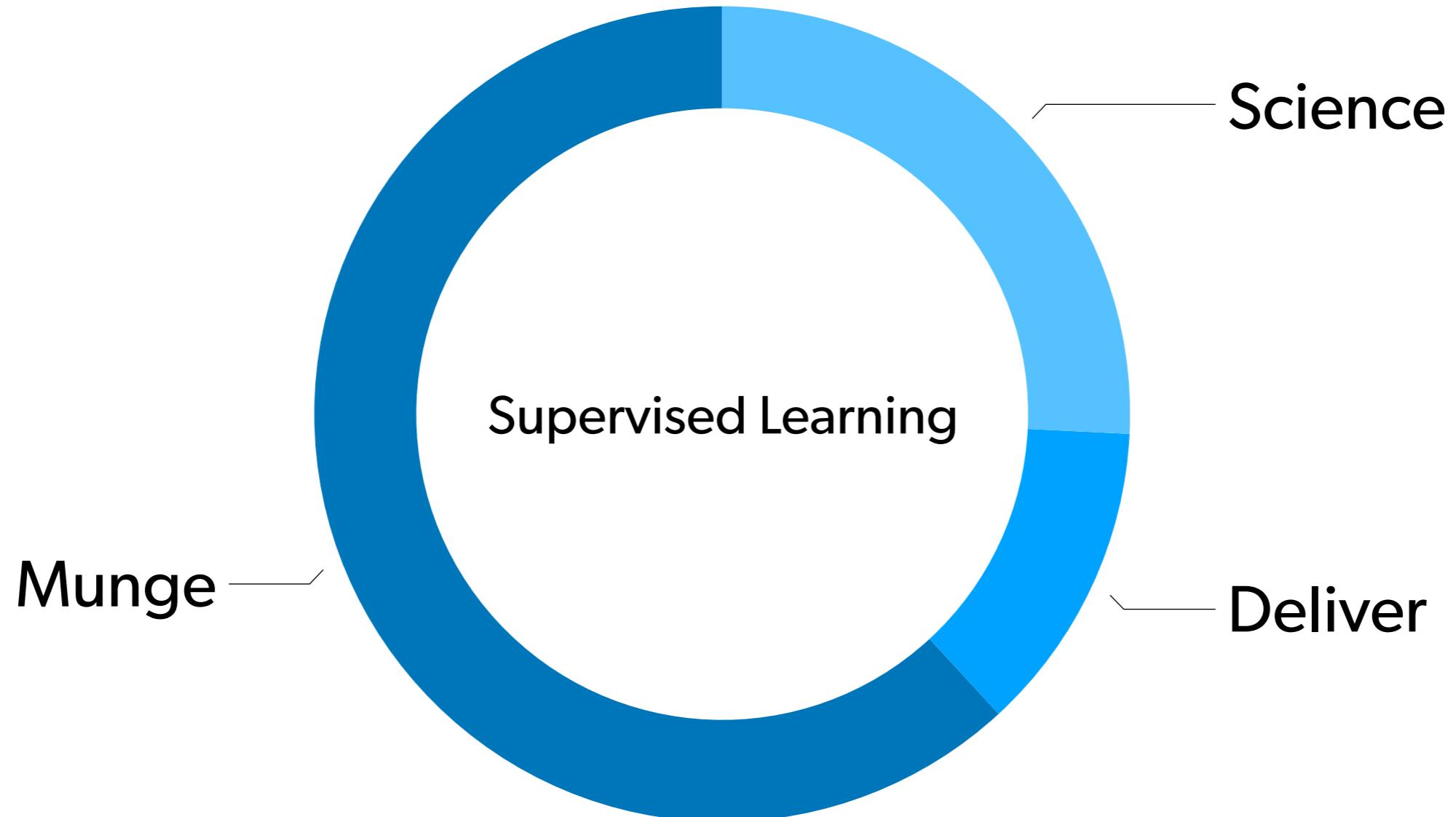
https://bit.ly/bus_number

Goal	Customer	Deliverables	Why Reproducible?
Exploratory Data Analysis	You	Intuition, knowledge, interactive widgets, reusable code	To help future you
Build, Deploy Models	Your Organization	Trained models, reliable metrics	Easy to update, easy to hand off
Publish or Share	Other (Data) Scientists	Paper, slides, code	Science without reproducibility is magic

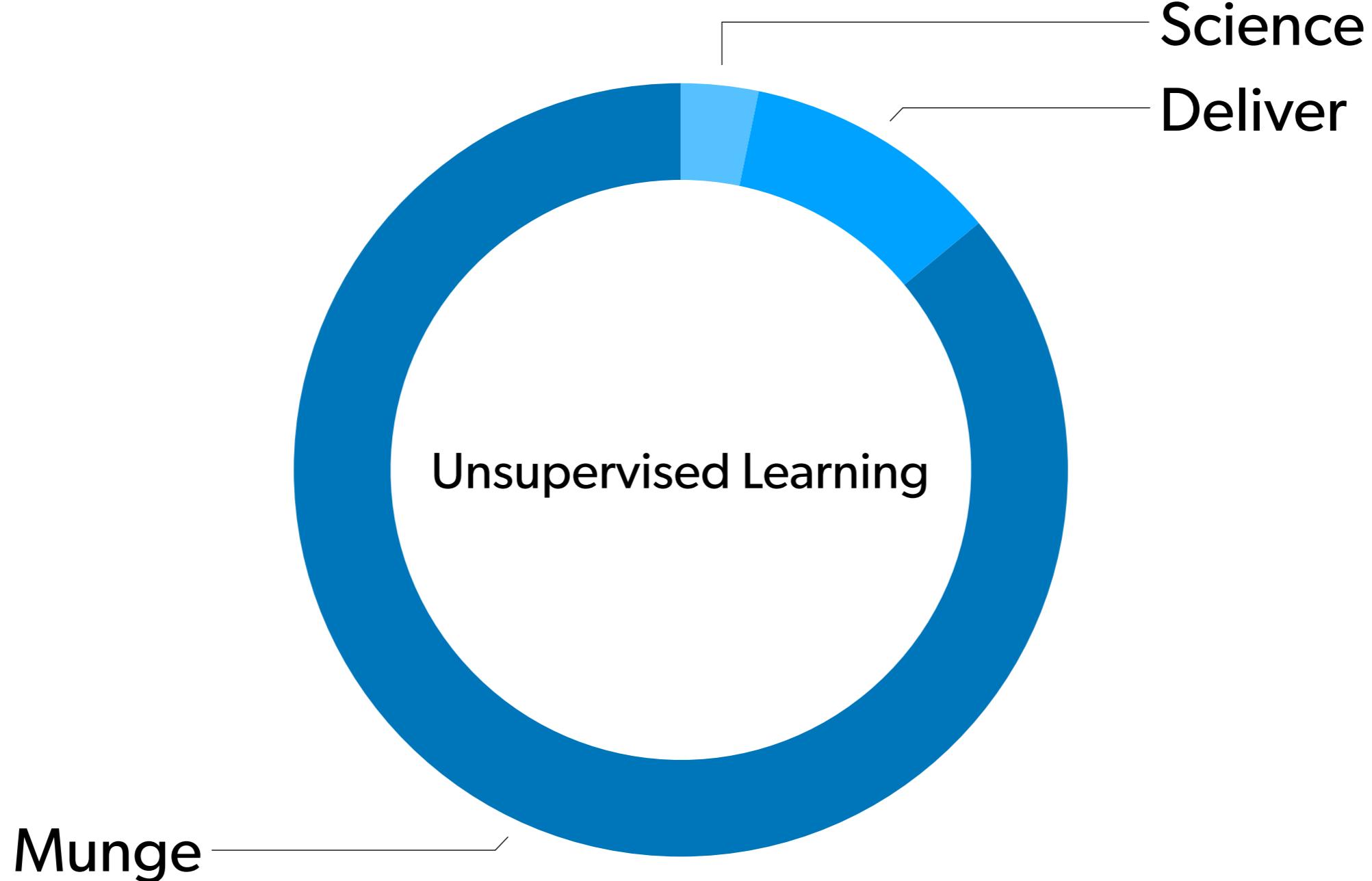
jupyter notebook

https://bit.ly/bus_number

How do you spend your “Data Science” time?



How do you spend your “Data Science” time?





Source: Youtube: Cårven Der Pümpkin: <https://www.youtube.com/watch?v=2Qj8PhxSnhg>

The Right Tools for the Job

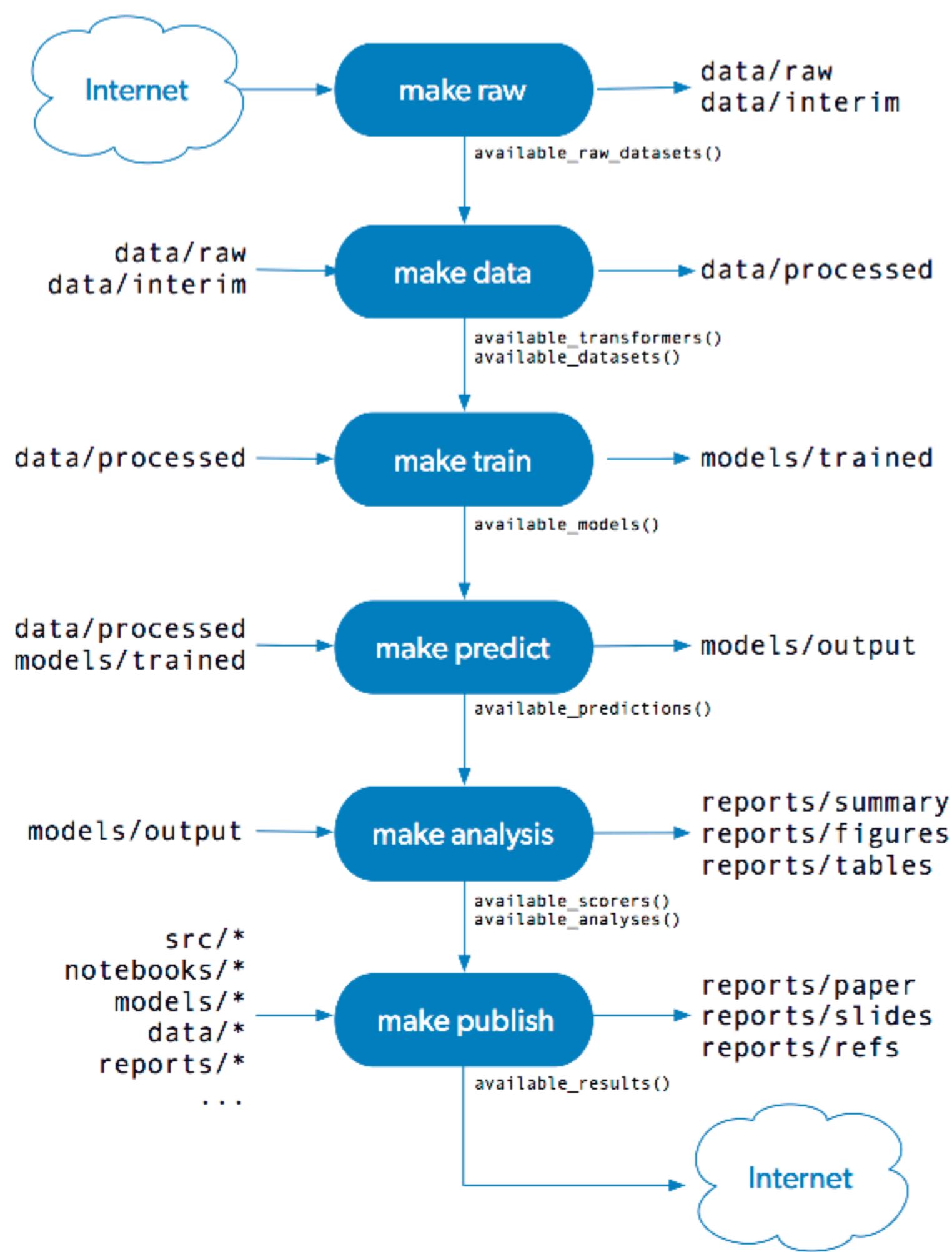
https://bit.ly/bus_number

The Right Tools for the Job

- Revision Control: **git** and (github/gitlab/bitbucket)
- Language: **python 3.6+**
- Package Manager and Virtual Environments: **conda**
- Frameworks: **scikit-learn, joblib**
- IDE: **jupyter notebook**
- Scripting: **make**
- Templates: **cookiecutter**

Data Science is a DAG

https://bit.ly/bus_number

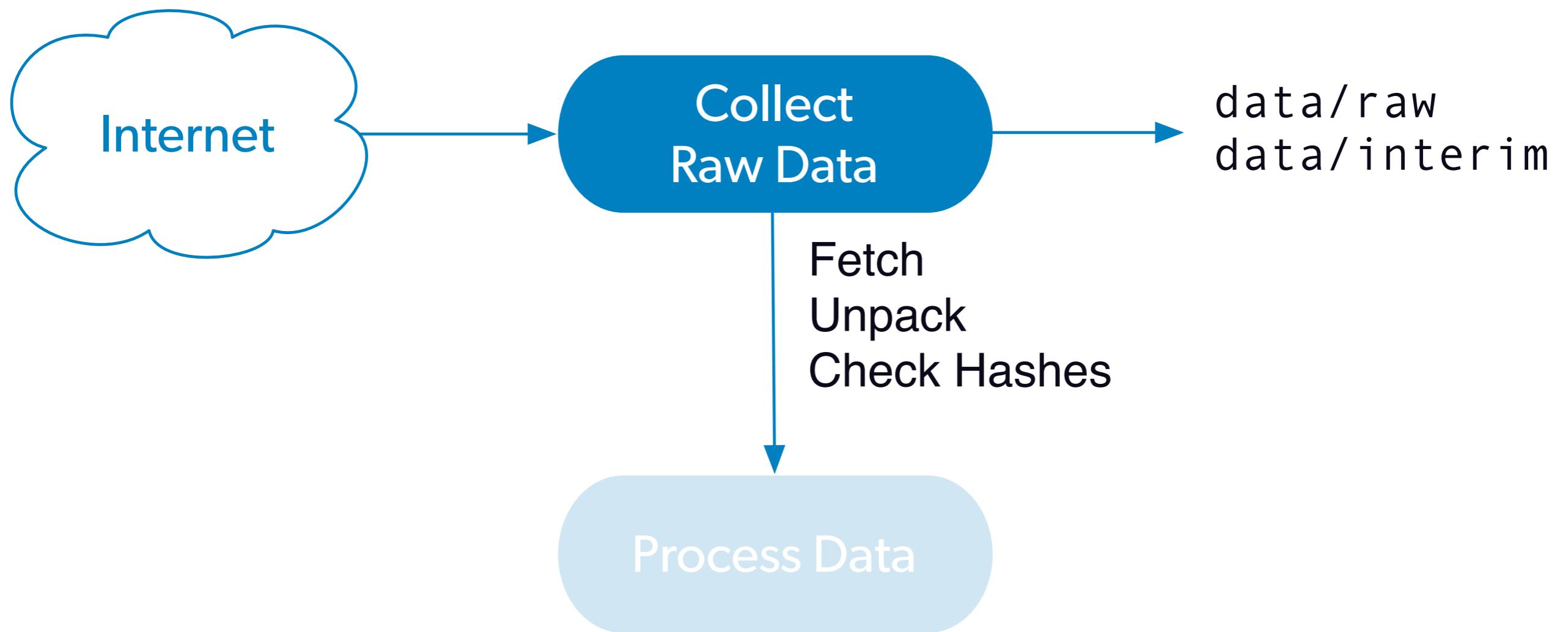


“Raw Data is Read-Only”

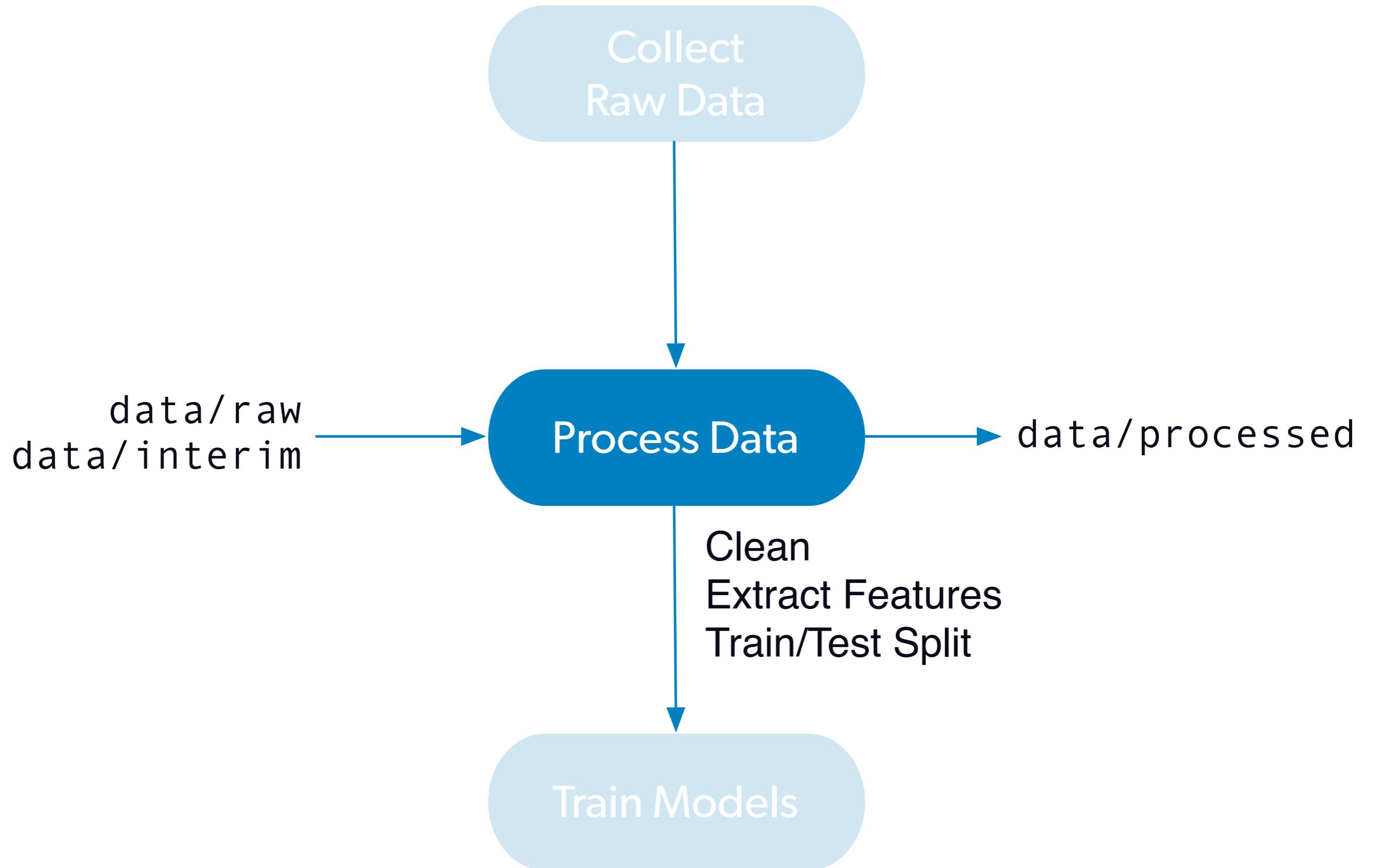


Sing it with me

https://bit.ly/bus_number

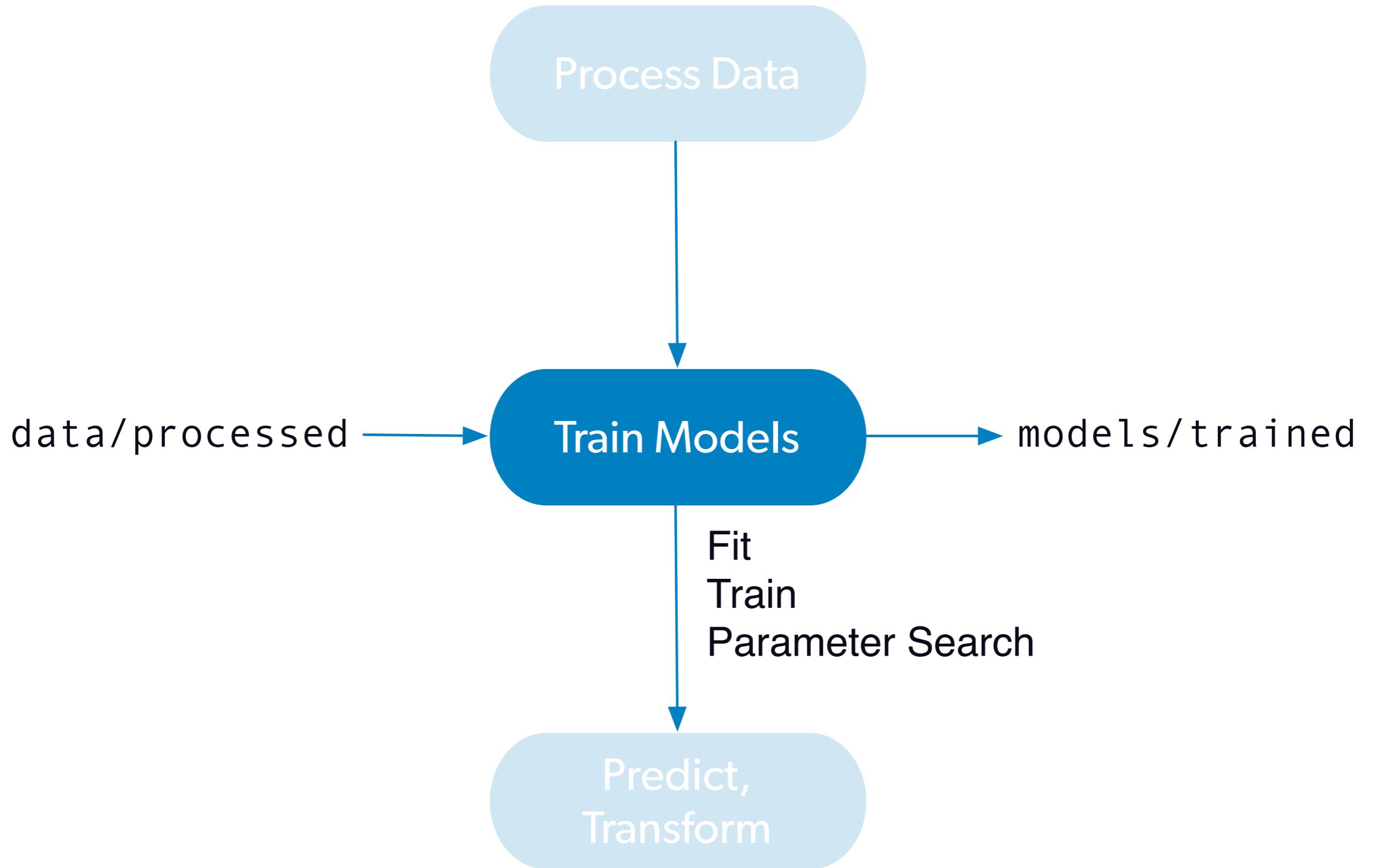


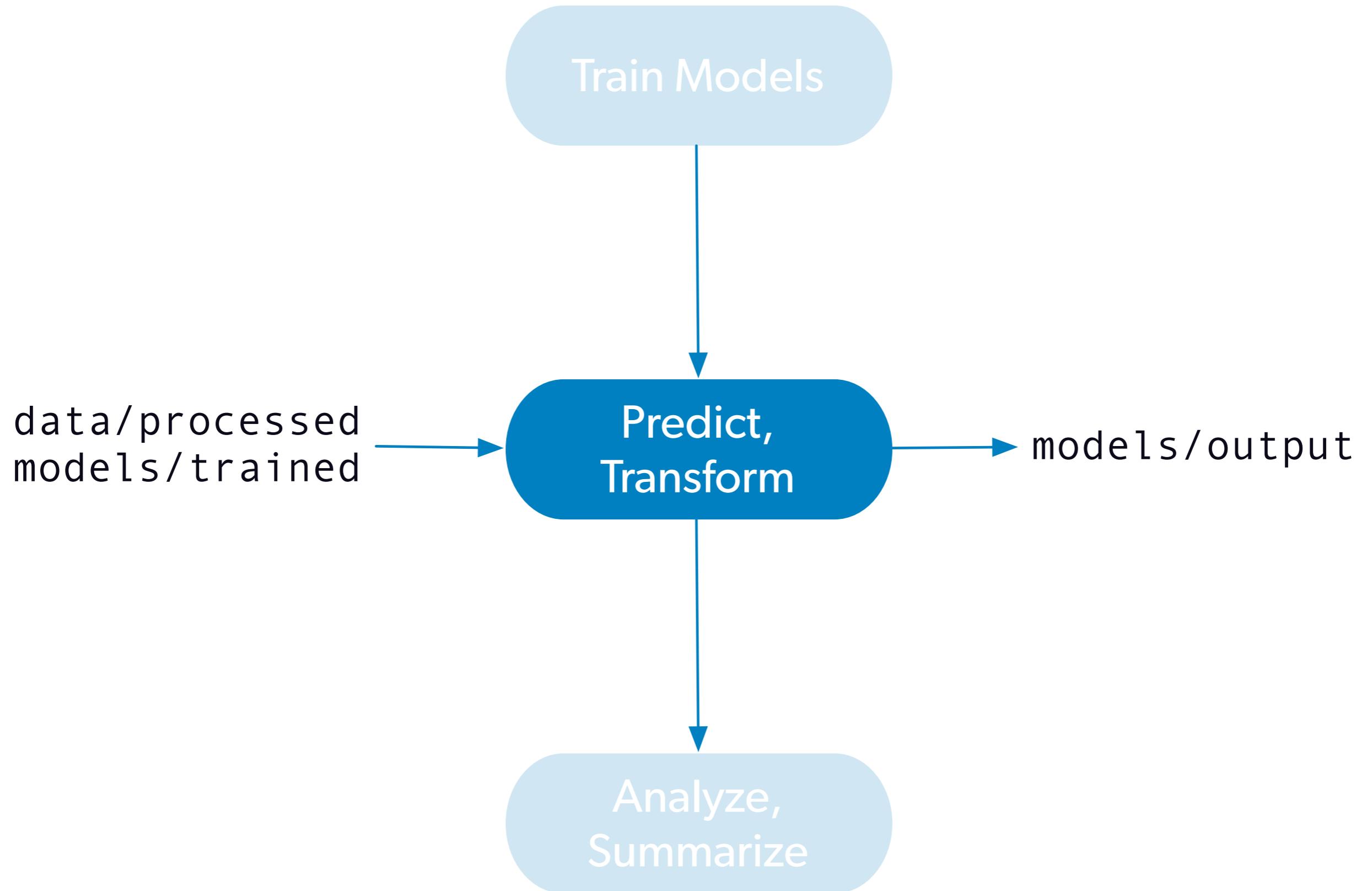
https://bit.ly/bus_number





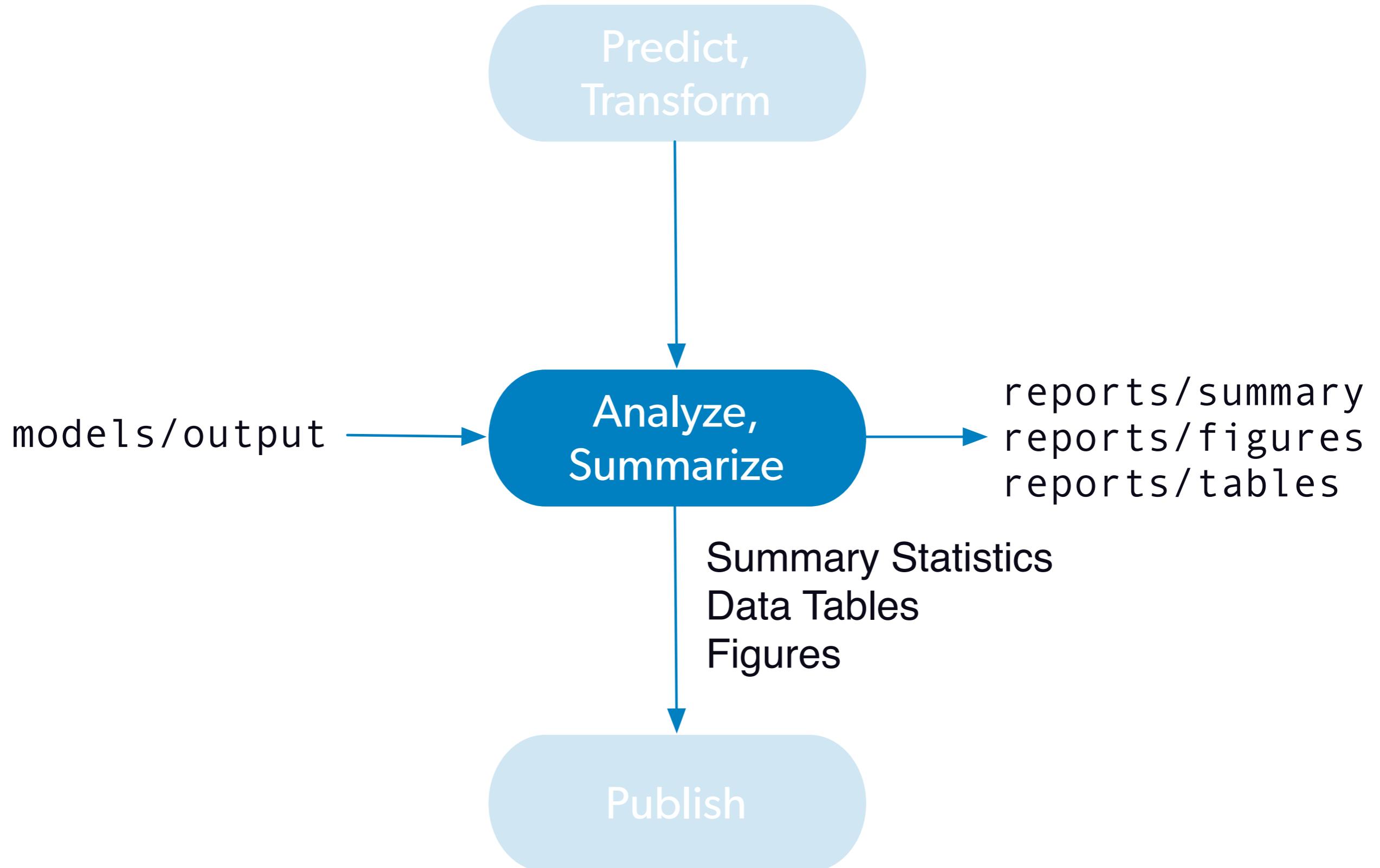
**We're gonna (data) science
the *@#! out of this**

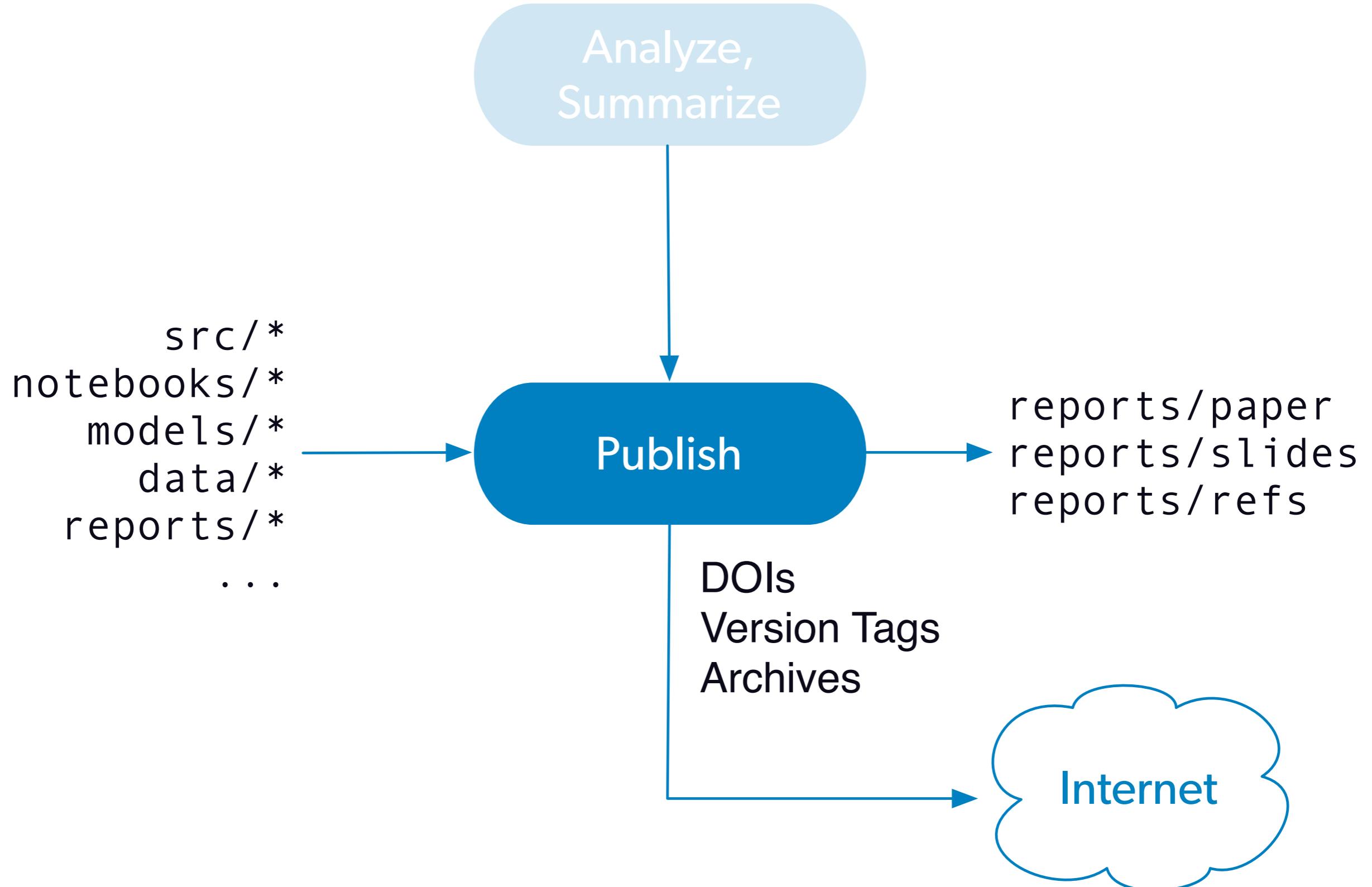




I'm only here
for the Pretty
Pictures







https://bit.ly/bus_number

What's your Bus Number?



Source: <https://www.imdb.com/title/tt0111257>

https://bit.ly/bus_number