

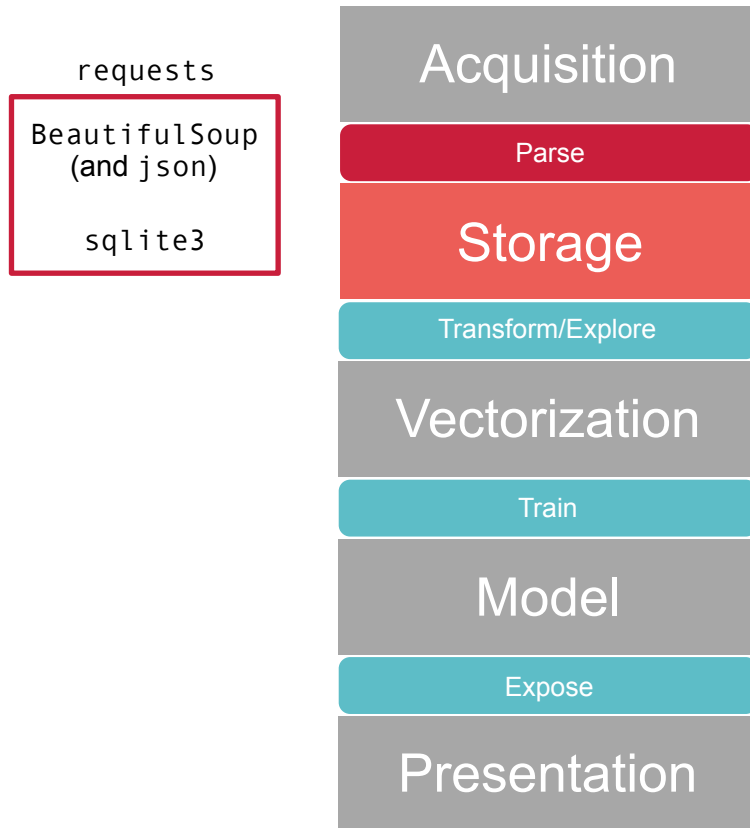
Lesson 4: Adding Structure— Parsing Data and Data Models

**For now, we are interested not in listings
themselves but the context around the
listing**

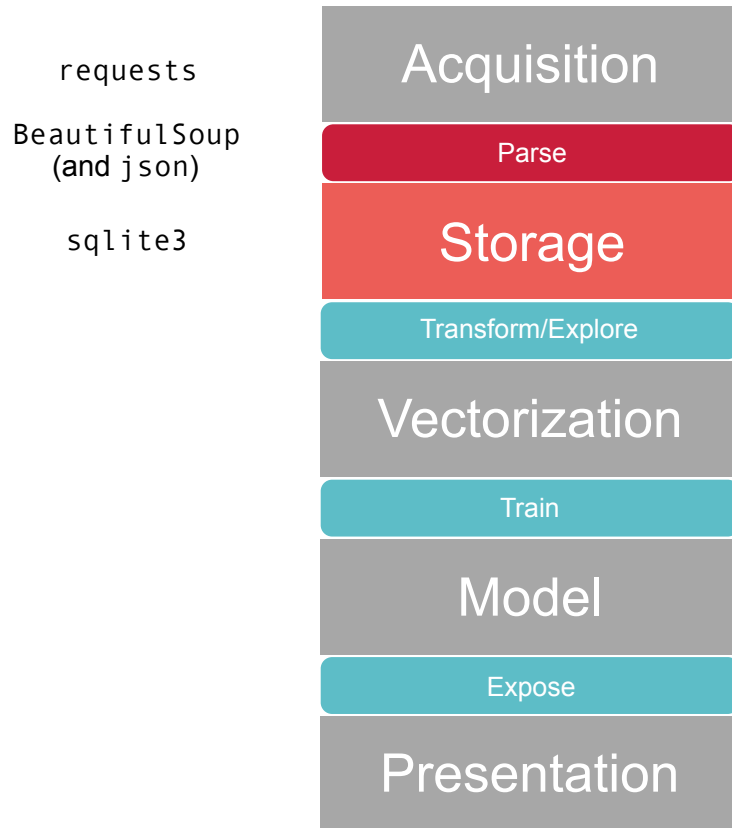
**Which neighborhoods are the most
explorable?**

Which listings have the most lively and active POI around them?

Process + Tools



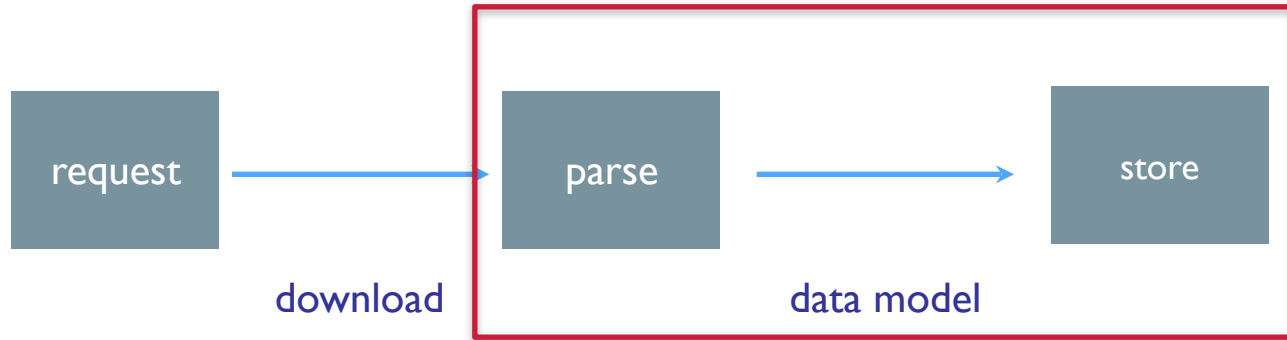
Process + Tools



ETL Pipeline



ETL Pipeline



Concepts

- Data Models
- Object Oriented Programming
- Schemas and Data Types
- Data Blending and Joining
- Types of Data (Time-series, Geo)
- The Shape of Data (wide vs. long)

Technologies

- Relational Databases (sqlite3)
- Object Relational Mappers (peewee)

Ideas

- Data Models
- Object Oriented Programming
- Schemas and Data Types
- Data Blending and Joining
- Types of Data (Time-series, Geo)
- The Shape of Data (wide vs. long)

Implementations

- Relational Databases (sqlite3)
- Object Relational Mappers (peewee)

ETL Pipeline



[illegible]

}

[illegible]

}

```

"email": "mailto:", "https://www.starbucks.com/menu/446232/cakesthe5157462730",
"email": "12063090",
"ip": "146.0.0.0/24", "https://www.starbucks.com/menu/446232/cakesthe5157462730",
"likes": {
  "count": 6,
  "groups": [],
  "summary": "0 likes"
},
"tagline": "true",
"text": "Better food, better coffee and better stuff. Why go to Starbucks down the street and waste your money?",
"count": 0
},
"type": "user",
"user":
  "first_name": "Marissa",
  "gender": "male",
  "id": "256438",
  "last_name": "Gorman",
  "gender": "male",
  "url": "https://1rx3.ack.net/Amg/user/",
  "url": "https://1rx3.ack.net/Amg/user/"
}

```

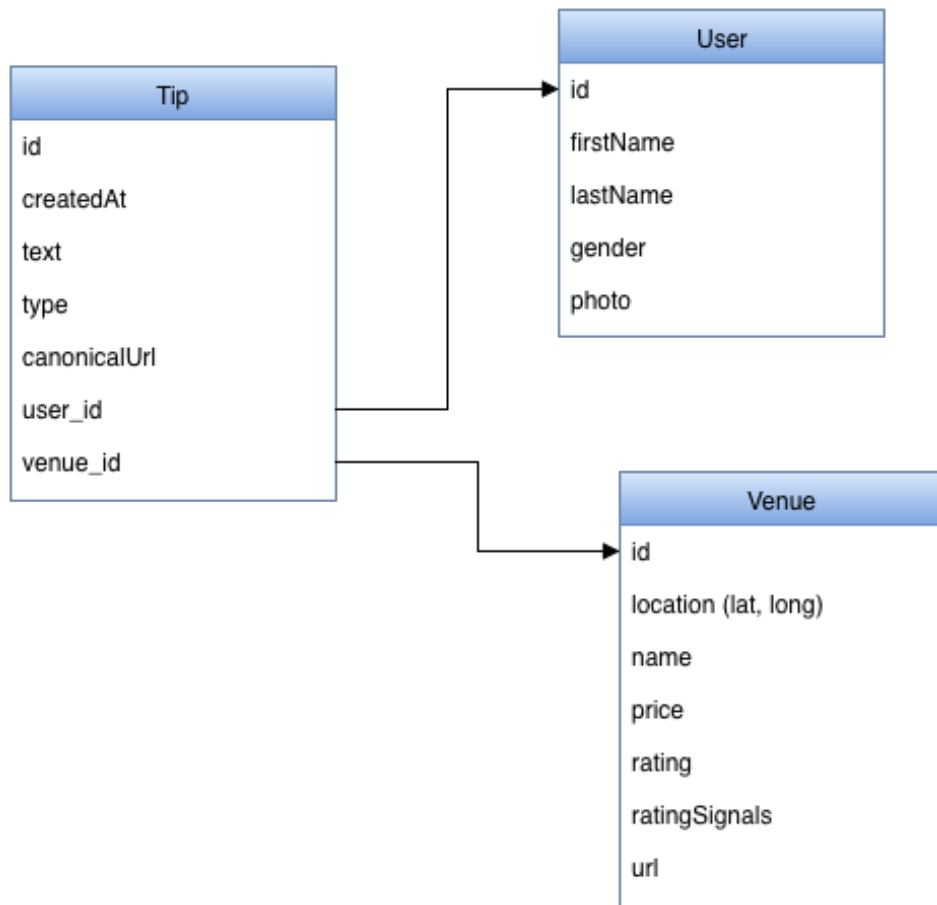
}

[illegible]

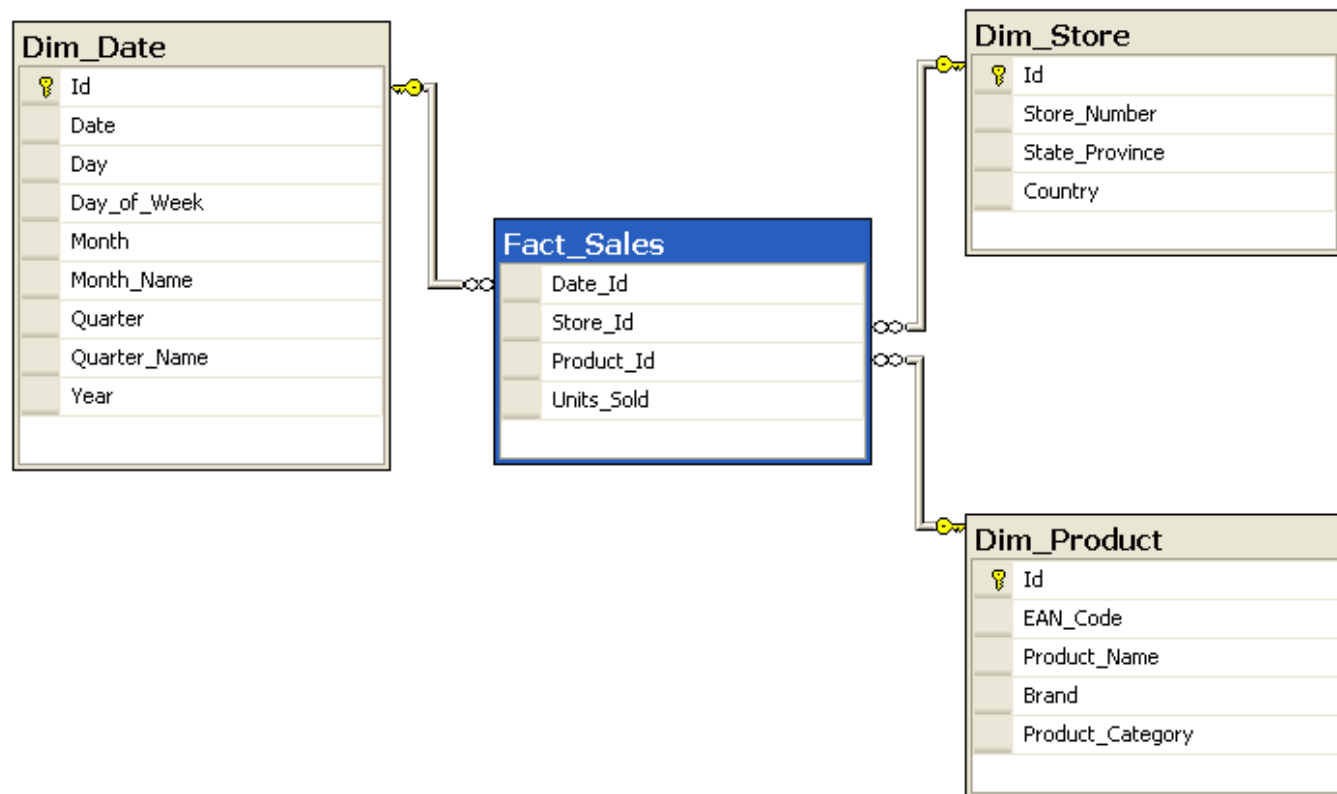
}

Thinking in Tables

Schema



Star Schema



JSON -> Data Tables

Functional Parsing

```
def parse_venue(data):  
    return {  
        'id': data['id'],  
        'name': data['name'],  
        'price': data['price']['tier'],  
        'rating': data['rating'],  
        'ratingSignals': data['ratingSignals'],  
        'url': data['url'],  
        'latitude' : data['location']['lat'],  
        'longitude': data['location']['lng']  
    }
```

Functional Parsing

```
def parse_user(data):  
    return {  
        'id': data['id'],  
        'firstName': data['firstName'],  
        'lastName': data['lastName'],  
        'gender': data['gender'],  
        'photo': data['photo']['prefix'] + 'original' + data['photo']['suffix']  
    }
```

Functional Parsing

```
def parse_tip(data, whitelist=['id', 'createdAt', 'text', 'type', 'canonicalUrl']):  
    user = parse_user(data['user'])  
    tip = { key: data[key] for key in whitelist }  
  
    return (tip, user)
```

Functional Parsing

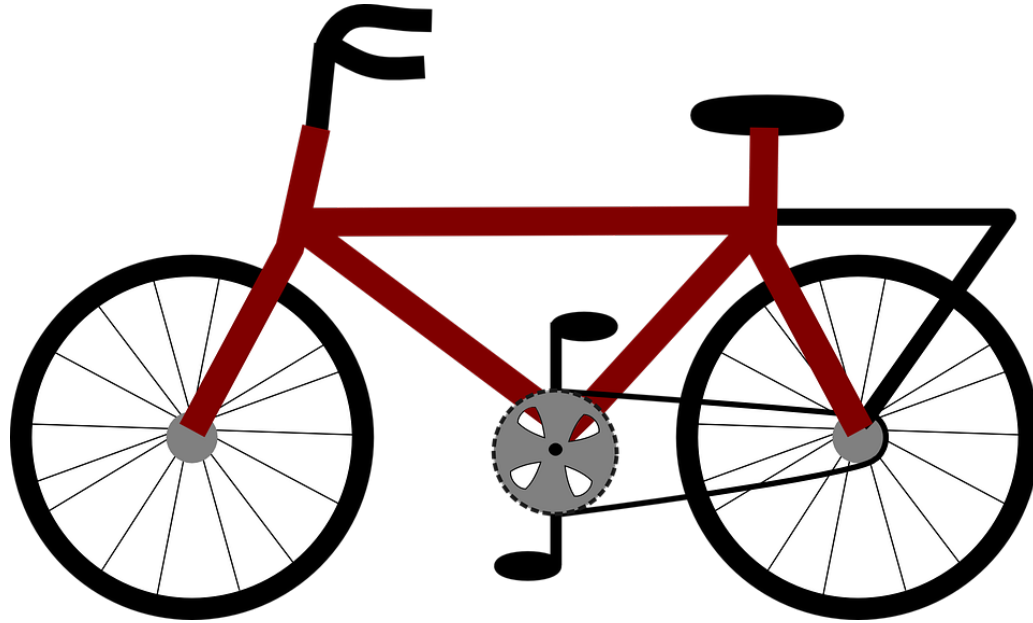
```
def parse_response(data):  
    tips = []  
    users = {}  
  
    venue = parse_venue(data['venue'])  
  
    for tip in data['tips']:  
        tip, user = parse_tip(tip)  
        tip['user_id'] = user['id']  
        tip['venue_id'] = venue['id']  
        tips.append(tip)  
  
        if user['id'] not in users:  
            users[user['id']] = user  
  
    return (venue, list(users.values()), tips)
```

Object Oriented Programming

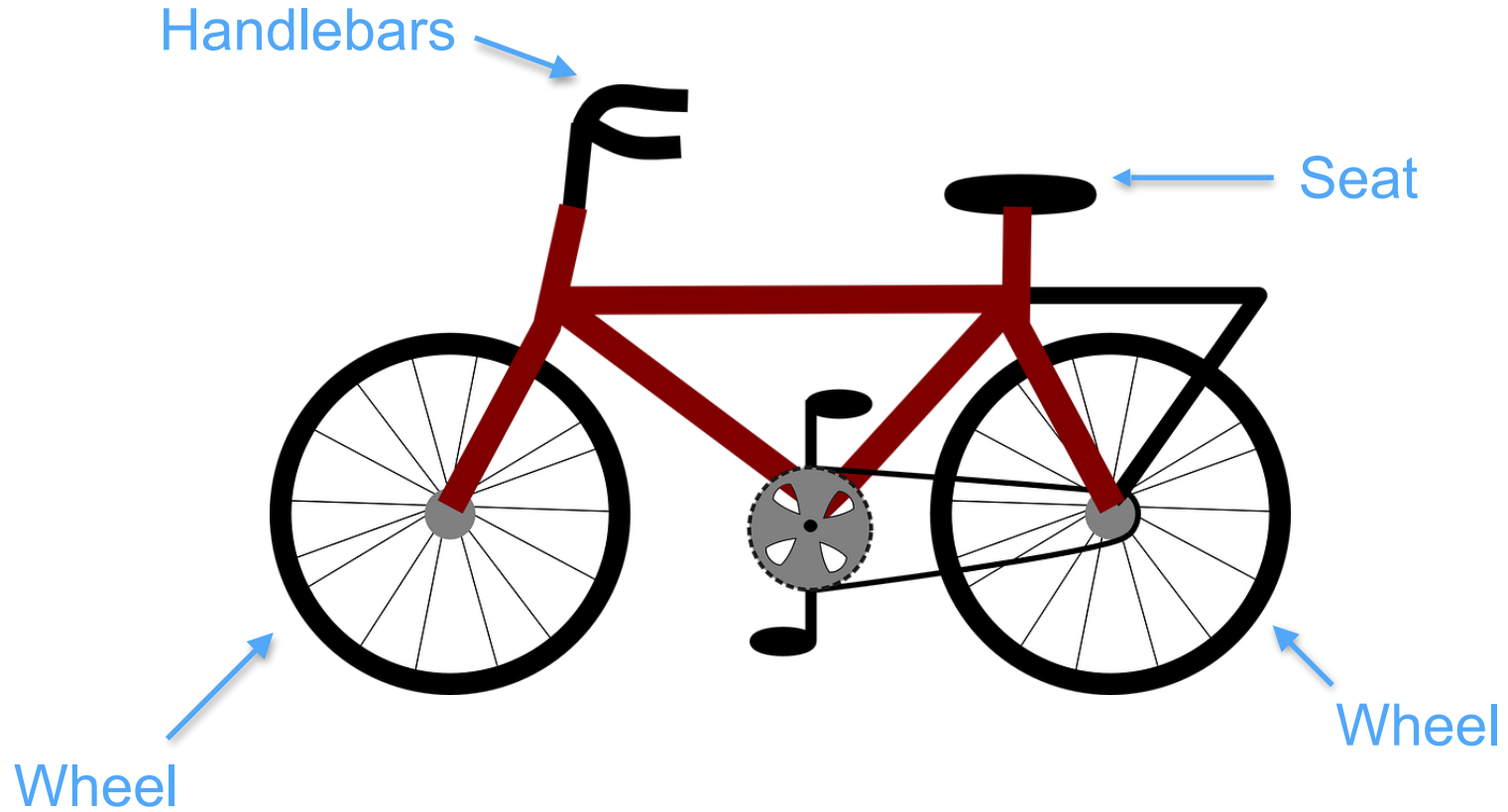
Object-oriented programming (OOP) is a programming paradigm based on the concept of "objects", which may contain data, in the form of fields, often known as attributes; and code, in the form of procedures, often known as methods.

- Wikipedia

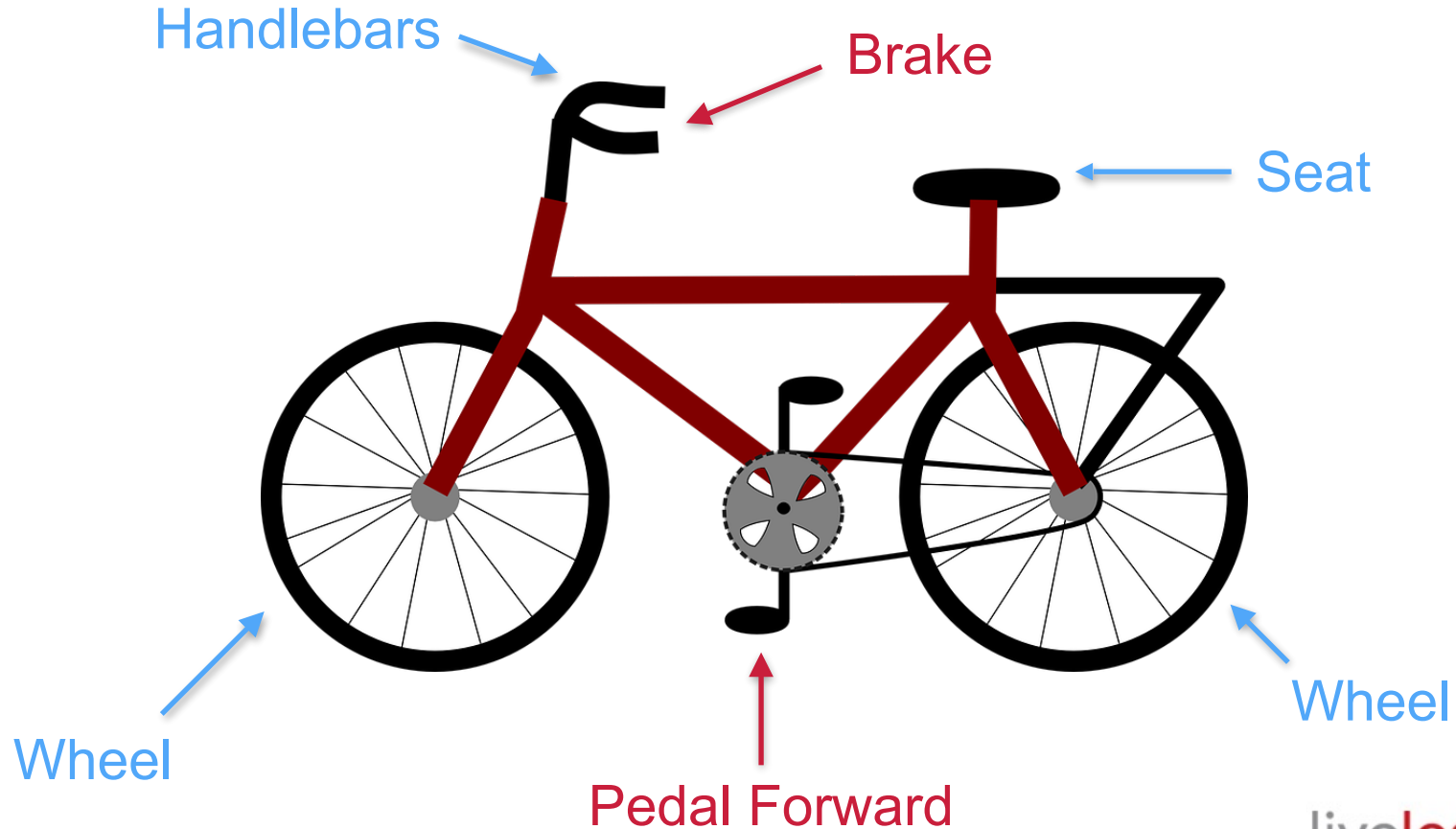
An Object Oriented Bicycle



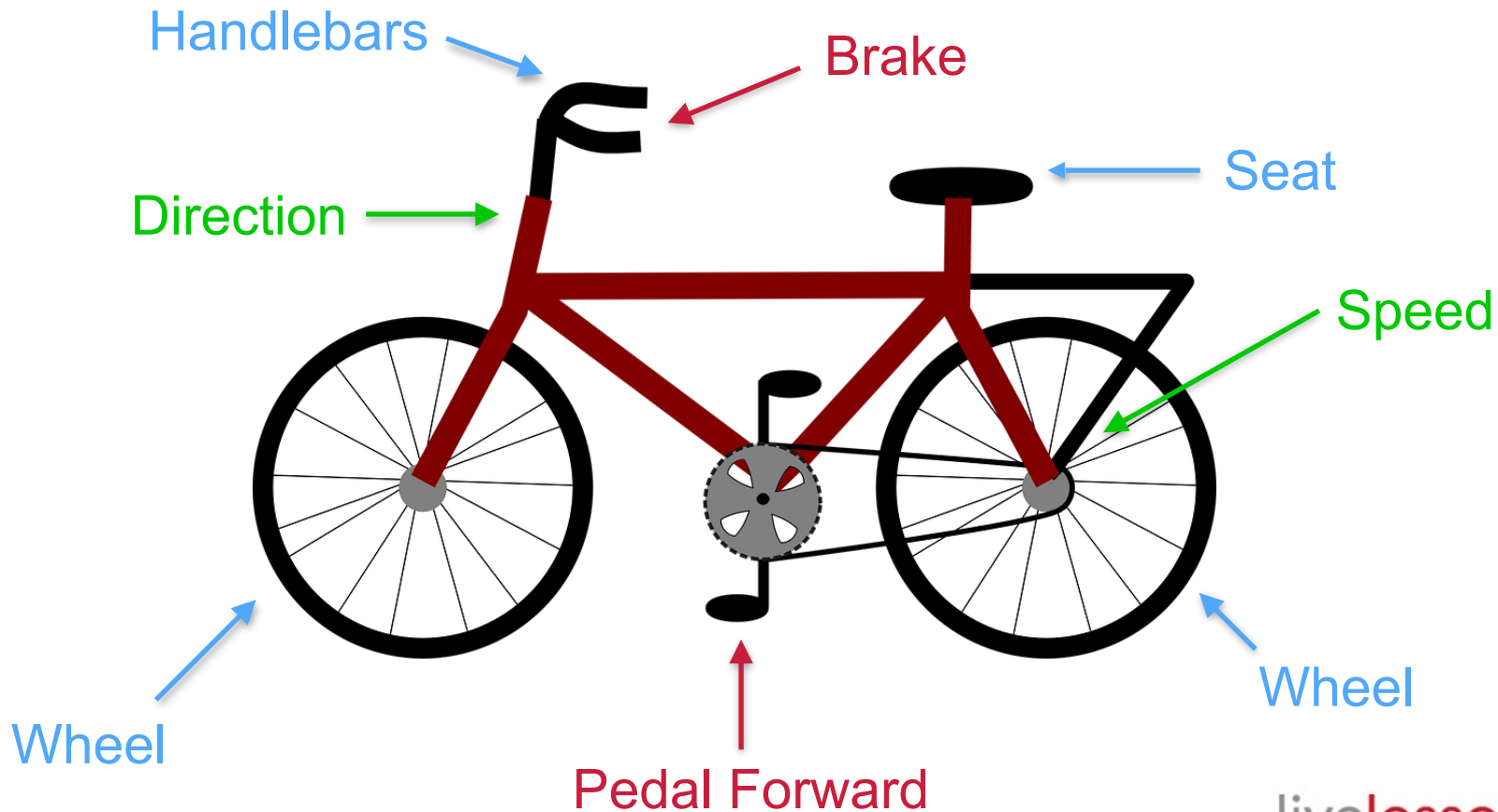
Attributes



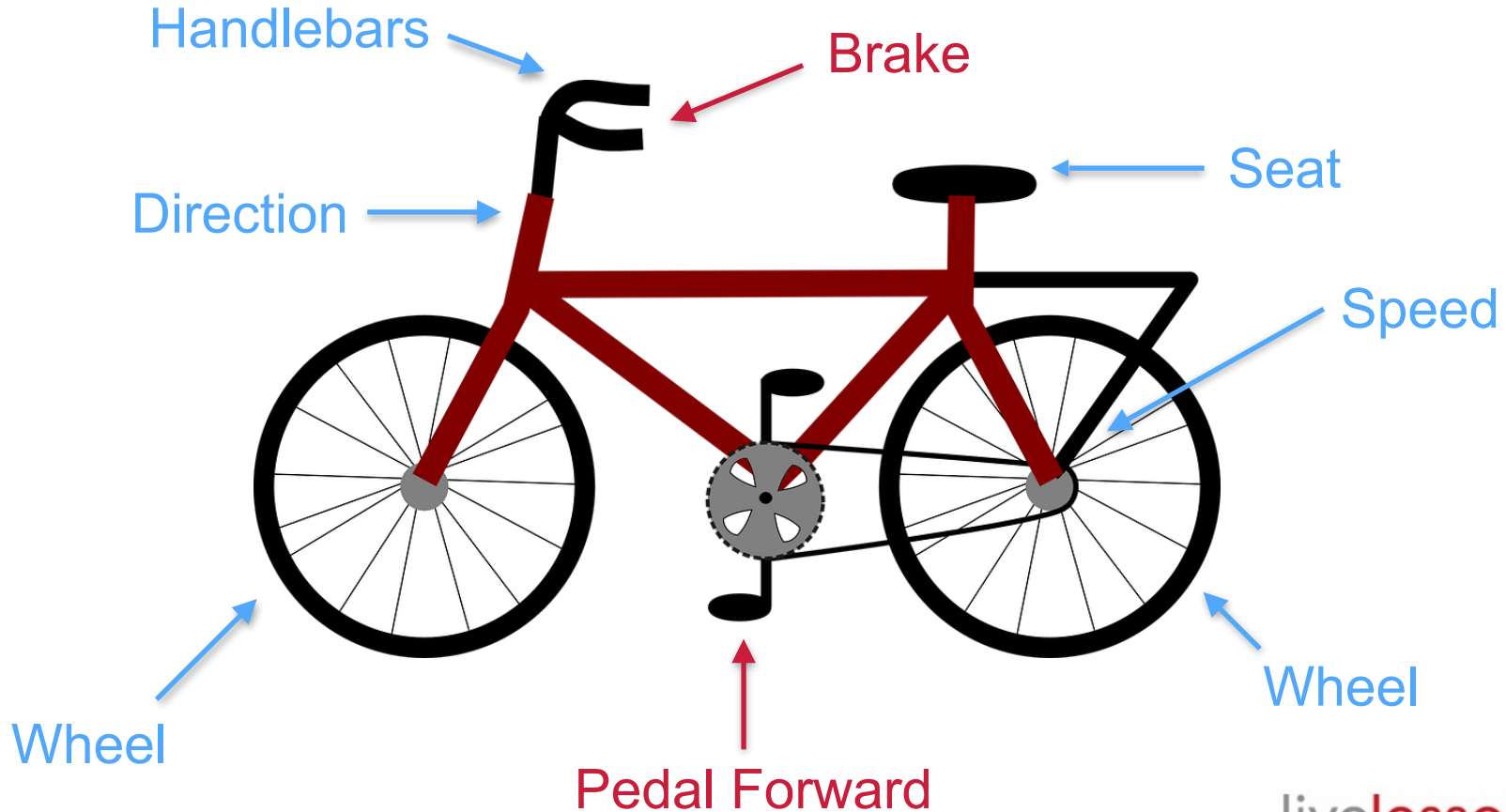
Methods



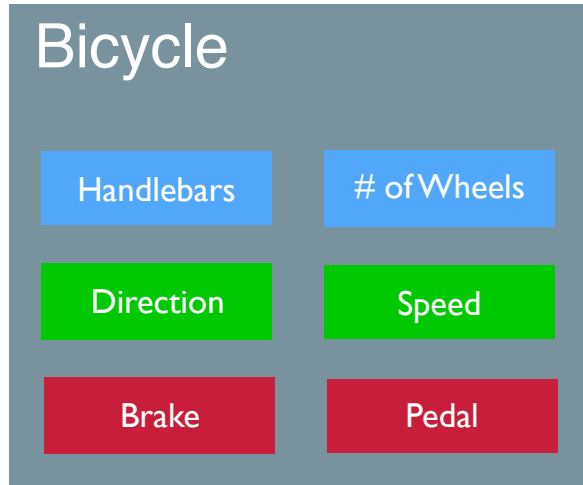
State



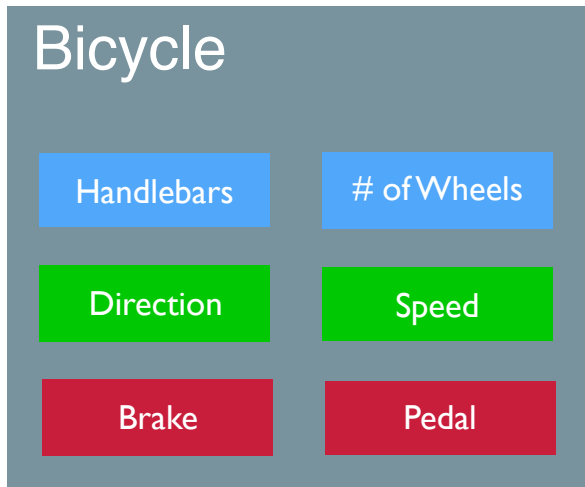
State (stored in Attributes)



Abstracted



Abstracted

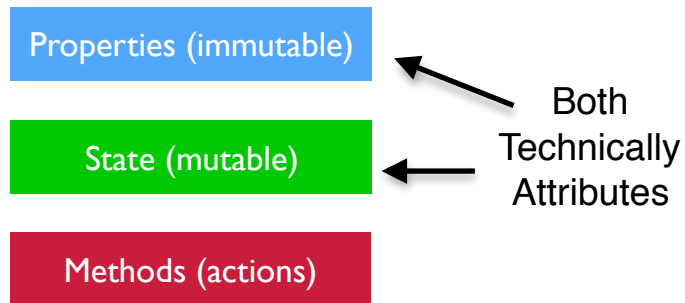
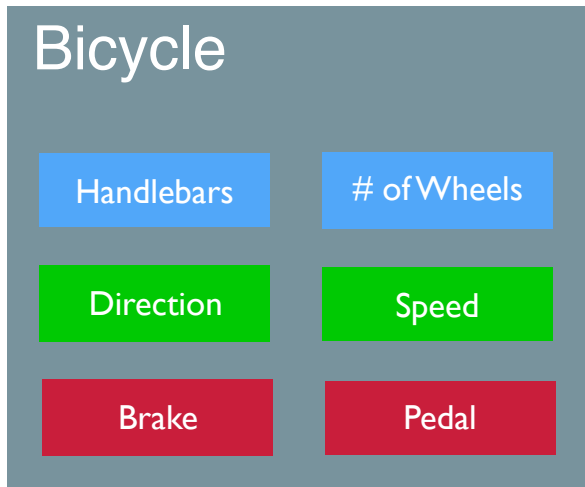


Properties (immutable)

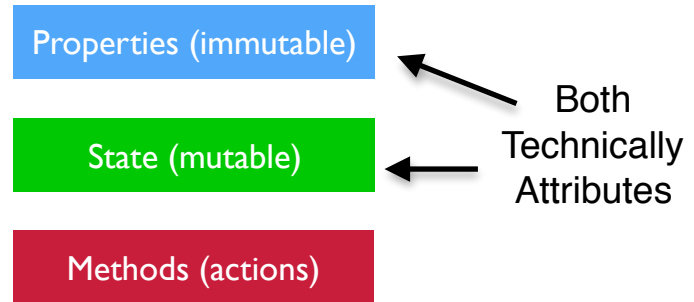
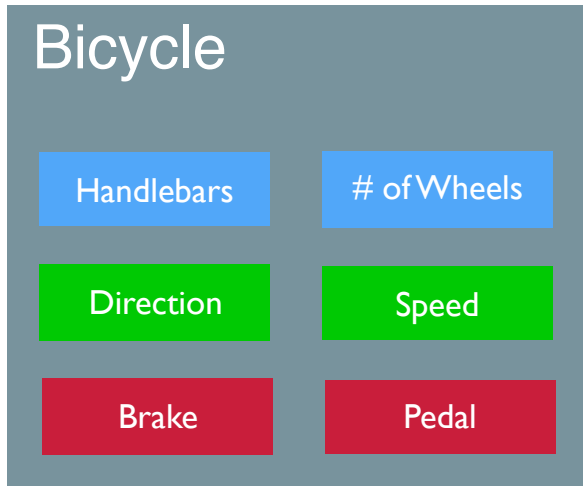
State (mutable)

Methods (actions)

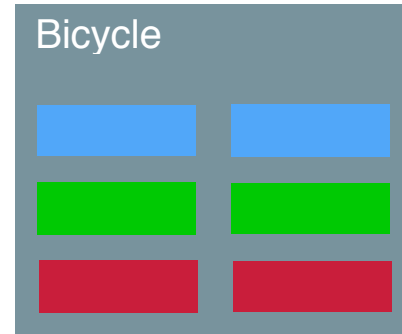
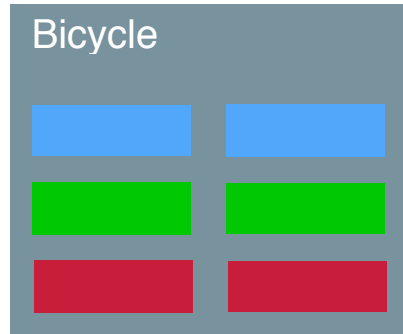
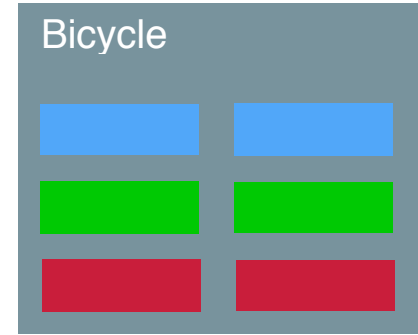
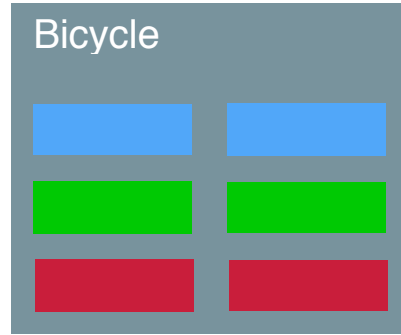
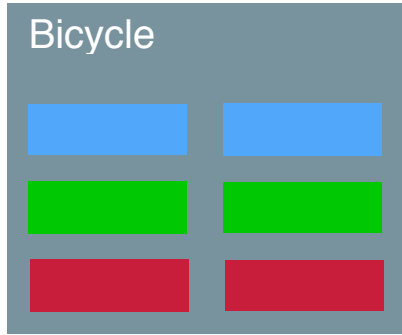
Abstracted



Class (think blueprint)



Objects and Instances



**In Python... everything is
(secretly) an object!**

**Lets talk about when we might benefit
from OOP principles by applying them to
our Foursquare Data**

Makes some things easier to program

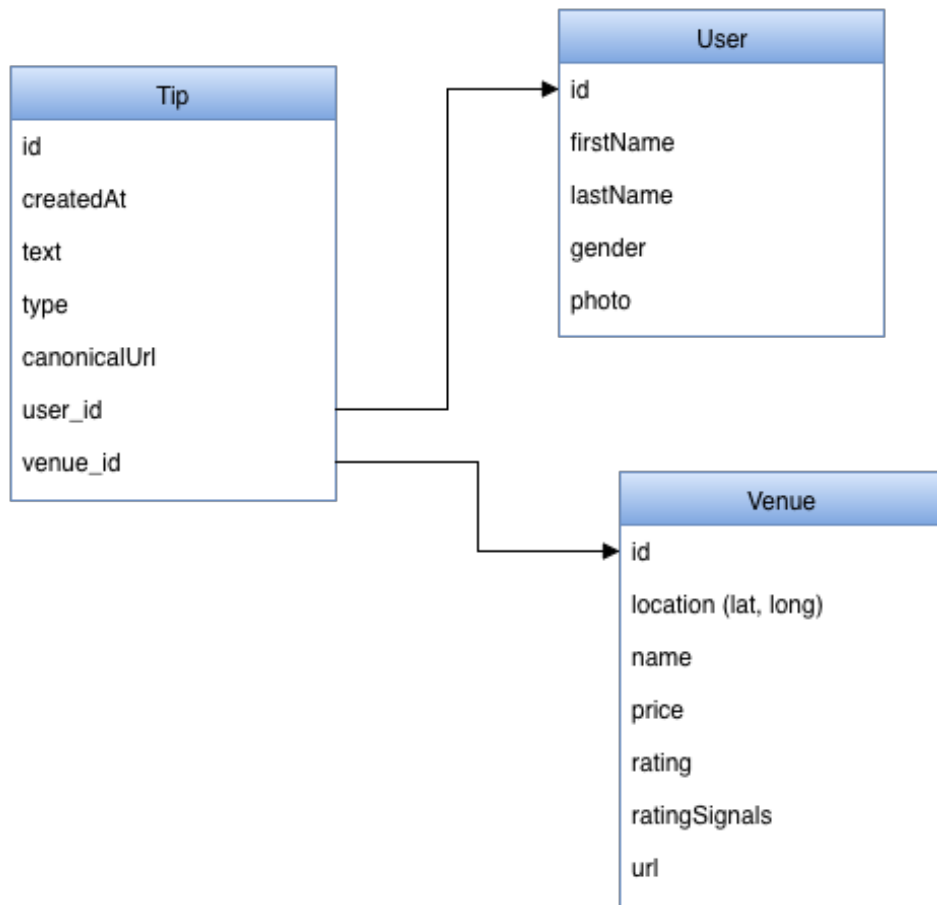
Makes other things harder to program

OOP allows you to combine **data and
functionality**

OOP allows you to combine attributes and methods

**How might we represent our
Foursquare data as **objects**?**

Schema



Keyword Arguments

```
class Venue:
    def __init__(self, **data):
        self.id = data['id']
        self.name = data['name']
        self.price = data['price']
        self.rating = data['rating']
        self.ratingSignals = data['ratingSignals']
        self.url = data['url']
        self.latitude = data['latitude']
        self.longitude = data['longitude']
```

```
class User:
    def __init__(self, **data):
        self.id = data['id']
        self.firstName = data['firstName']
        self.lastName = data['lastName']
        self.gender = data['gender']
        self.photo = data['photo']
```

```
class Tip:
    def __init__(self, **data):
        self.id = data['id']
        self.createdAt = data['createdAt']
        self.text = data['text']
        self.canonicalUrl = data['canonicalUrl']
```

**Lets talk about when we might benefit
from OOP principles by applying them to
our Foursquare Data**

**Live Coding (in sublime): Just
Data, lets add methods
(fsmodels.py)**

class Tip:

def __init__(self, **data):

self.id = data['id']

self.createdAt = data['createdAt']

self.text = data['text']

self.canonicalUrl = data['canonicalUrl']

relations

self.user = data['user']

self.location = data['venue']

set initial values

self.upvotes = 0

self.downvotes = 0

def upvote(self):

self.upvotes += 1

def downvote(self):

self.downvote += 1

jessica = User(**u)

cafe = Venue(**v)

tip = Tip(**t, user = jessica, venue = cafe)

Inheritance

```
class Tip:
    def __init__(self, **data):
        self.id = data['id']
        self.createdAt = data['createdAt']
        self.text = data['text']
        self.canonicalUrl = data['canonicalUrl']

        # relations
        self.user = data['user']
        self.location = data['venue']

        # set initial values
        self.upvotes = 0
        self.downvotes = 0

    def upvote(self):
        self.upvotes += 1

    def downvote(self):
        self.downvote += 1
```

Inheritance

```
class Checkin:
    def __init__(self, **data):
        self.createdAt = time.time()

        # relations
        self.user = data['user']
        self.location = data['venue']
```

Inheritance

```
class Tip:
    def __init__(self, **data):
        self.id = data['id']
        self.createdAt = data['createdAt']
        self.text = data['text']
        self.canonicalUrl = data['canonicalUrl']

        # relations
        self.user = data['user']
        self.location = data['venue']

        # set initial values
        self.upvotes = 0
        self.downvotes = 0

    def upvote(self):
        self.upvotes += 1

    def downvote(self):
        self.downvote += 1
```

```
class Checkin:
    def __init__(self, **data):
        self.createdAt = time.time()

        # relations
        self.user = data['user']
        self.location = data['venue']
```

Inheritance

```
class Tip:
    def __init__(self, **data):
        self.id = data['id']
        self.createdAt = data['createdAt']
        self.text = data['text']
        self.canonicalUrl = data['canonicalUrl']

        # relations
        self.user = data['user']
        self.location = data['venue']

        # set initial values
        self.upvotes = 0
        self.downvotes = 0

    def upvote(self):
        self.upvotes += 1

    def downvote(self):
        self.downvote += 1
```

```
class Checkin:
    def __init__(self, **data):
        self.createdAt = time.time()

        # relations
        self.user = data['user']
        self.location = data['venue']
```

**OOP allows you to group/collect/arrange
code in extensible patterns that reduce
duplication**

**And less code means less chances for error
(usually)**

Use Cases

- OO Abstractions Naturally Fit Domain Abstractions (i.e. Video Games and GUI)
- Benefits of OOP Suit Problem Specifications (i.e. `scikit-learn`)

Use Cases

- OO Abstractions Naturally Fit Domain Abstractions (i.e. Video Games and GUI)
- **Benefits of OOP Suit Problem Specifications (i.e. scikit-learn)**

How can we design a system/program that is consistent, interpretable, and extensible in the face of updating requirements?

e.g. Designing a Machine Learning Library

Benefits

- Encapsulation and Modularization
- Inheritance and Extensibility (DRY)
- Composition and Separation of Concerns

When to Use

- Application Programming Interfaces
- Multiple Independent (interacting) Agents
- Concurrent Components (through message passing)
- Long Lived Stateful Objects

When to Use

- Data Structures, Libraries, Frameworks, etc.
- Simulations (and Games)
- Distributed Systems*
- Graphics

* Implementation Dependent (relies on asynchronous messaging)

When to Use

- Custom Modules and Packages
- Poker and Blackjack (i.e. expected returns of strategies)
- Computer Multitasking (and [Operating System Interfaces](#))
- GUI windows

When to Use

- Custom Modules and Packages
- Poker and Blackjack (i.e. expected returns of strategies)
- Computer Multitasking (and [Operating System Interfaces](#))
- GUI windows

When Not to Use

- Small Scripts and “Glue” code
- Anything Procedural (think recipes and algorithms)
- Anything Functional (can you write it as an equation?)
- Data flow and Stream Processing (data pipelines)

https://en.wikipedia.org/wiki/Object-oriented_programming#Criticism

**At the end of the day it is mainly
personal preference**

**Remember... quick iteration is
key to data science!**

**Avoid tight coupling and be
careful with mutating data**

**And you can happily mix paradigms
(scikit-learn **does**)!**

Object Oriented Programming

The big idea is "messaging" -- that is what the kernal of Smalltalk/Squeak is all about (and it's something that was never quite completed in our Xerox PARC phase)... The key in making great and growable systems is much more to design how its modules communicate rather than what their internal properties and behaviors should be.

- Alan Kay

Object Oriented Programming

Think of the internet -- to live, it (a) has to allow many different kinds of ideas and realizations that are beyond any single standard and (b) to allow varying degrees of safe interoperability between these ideas.

- Alan Kay