

Machine Learning in Python with No Strings Attached

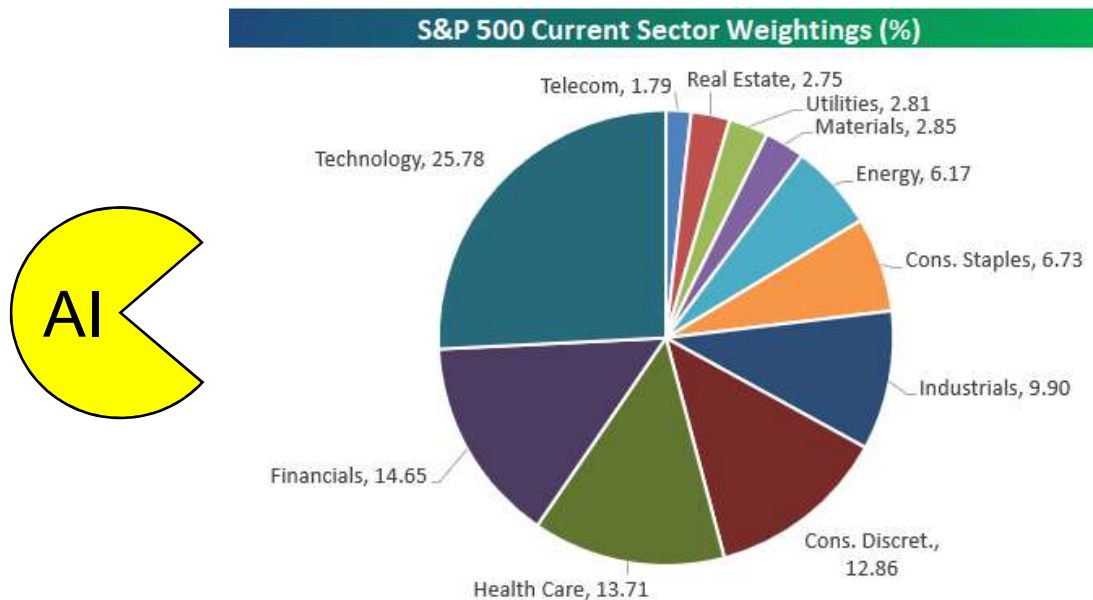
Guillaume Baudart, Martin Hirzel, Kiran Kate,
Louis Mandel, and Avi Shinnar

IBM Research, USA

6/22/2019 talk at MAPL Workshop

Why PL for AI?

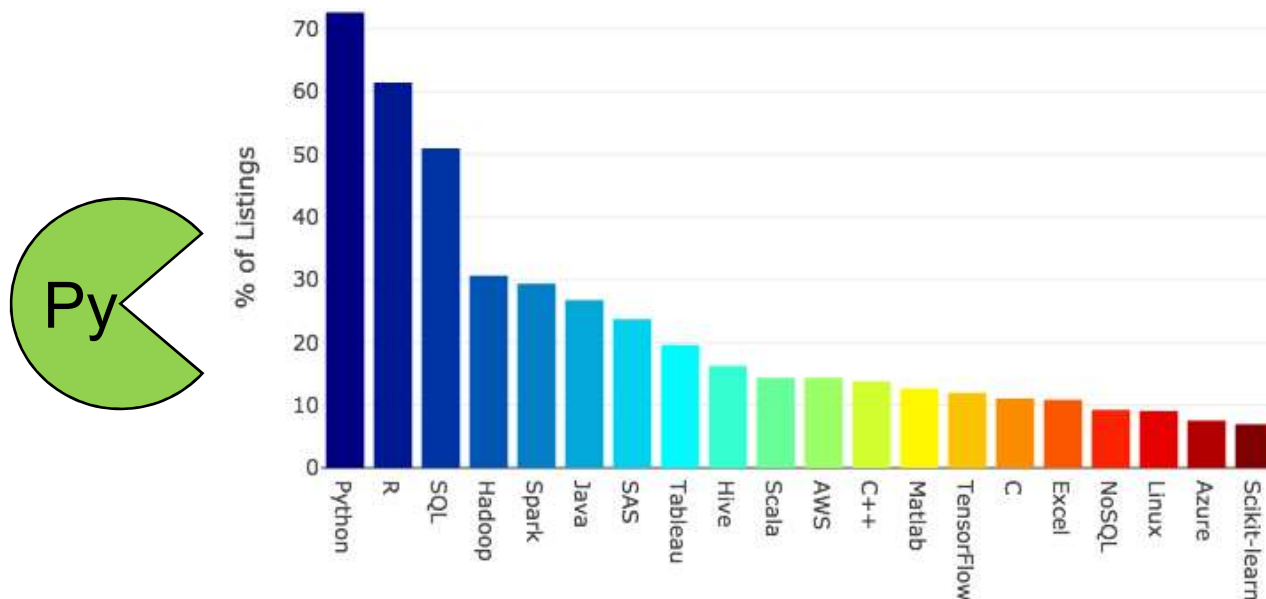
- IBM customers want AI
 - Large non-tech companies
 - Not enough skilled data scientists
- PL can lower the required skill level
 - Consistency, conciseness, error checking



Source: <https://seekingalpha.com/article/4172093-s-and-p-500-sector-weightings-tech-nears-26-percent>

Why Python?

- Good PL design: consistency, conciseness, ...
- Many libraries: Keras, Sklearn, Spark, ...
- Jupyter notebooks
- Popularity
 - Top 20 technology skills in data science job listings



Source: <https://towardsdatascience.com/the-most-in-demand-skills-for-data-scientists-4a4a8db896db>

Why Embedded Languages?

Build computational graph

- Host language: Python
- Guest language: Keras, Sklearn, Spark SQL, ...



Lazy evaluation

Evaluate

- Train vs. predict
- Run on GPU or cluster
- Save for later reuse
- Automatic differentiation
- Back-propagation
- Optimization
- Monte-Carlo sampling

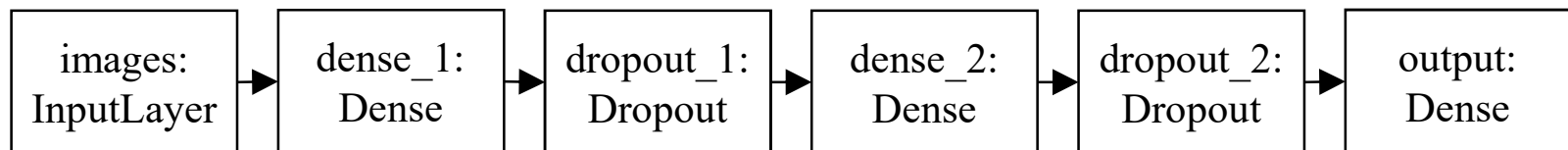
Just a library *–but–* also a language.

Keras Example

```
In [2]: import keras                                     1
        from keras.models import Model                  2
        from keras.layers import Dense, Dropout, Input 3
                                                4
        images = Input(shape=(N_PIXELS,), name='images') 5
        hidden1 = Dense(512, activation='relu')(images) 6
        dropped1 = Dropout(0.2)(hidden1)                7
        hidden2 = Dense(512, activation='relu')(dropped1) 8
        dropped2 = Dropout(0.2)(hidden2)                9
        output = Dense(N_LABELS, activation='softmax',   10
                        name='output')(dropped2)         11
                                                12
        model = Model(inputs=images, outputs=output)     13
```

```
In [3]: keras_utils.plot_model(model)
```

Out[3]:



```
In [5]: model.compile(loss='categorical_crossentropy', 1
                        optimizer='RMSprop',            2
                        metrics=['accuracy'])           3
        model.fit({'images': x_train}, {'output': y_train}, 4
                  batch_size=128, epochs=3,             5
                  validation_data=(x_test, y_test))      6
```

Outline

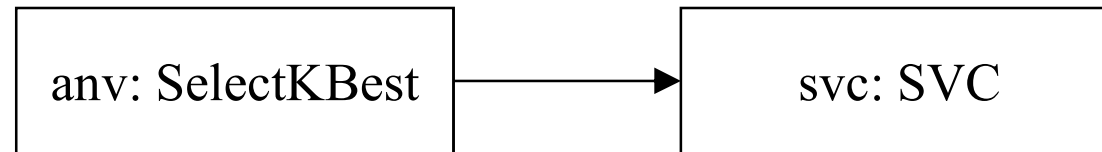
Strings Considered
Harmful

Name Reflection:
Keras.na

Reinterpreted
Python: YAPS

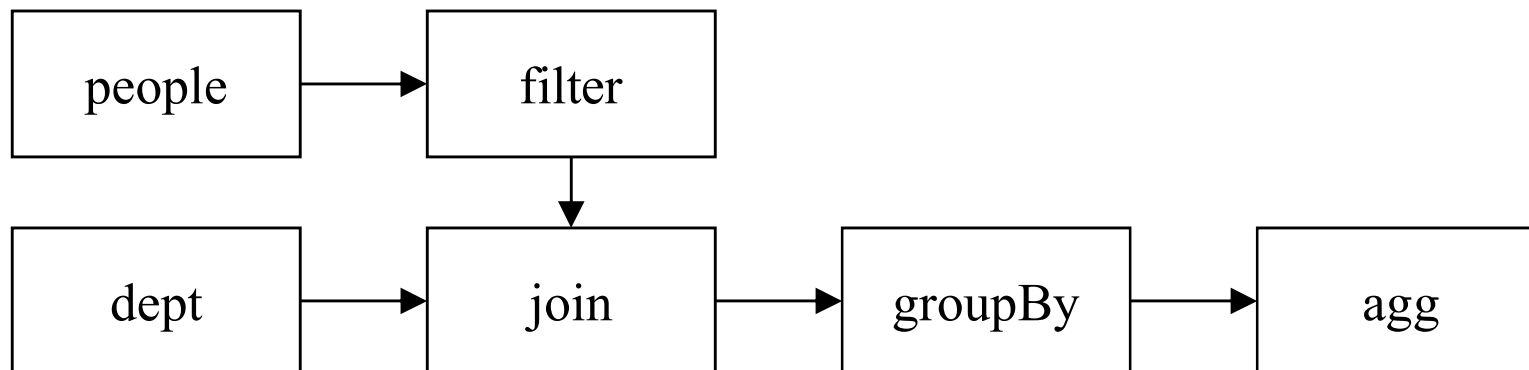
Sklearn Example

```
1 filter = SelectKBest(f_regression, k=5)
2 clf = svm.SVC(kernel='linear')
3 anv_svm = Pipeline([('anv', filter), ('svc', clf)])
4 anv_svm.set_params(anv__k=10, svc__C=.1).fit(X, y)
```



Spark SQL Example

```
1 people = sqlContext.read.parquet("people.parquet")
2 dept = sqlContext.read.parquet("dept.parquet")
3 people.filter(people.age > 30) \
4     .join(dept, people.deptId == dept.id) \
5     .groupBy(dept.name, "gender") \
6     .agg({"salary": "avg", "age": "max"})
```



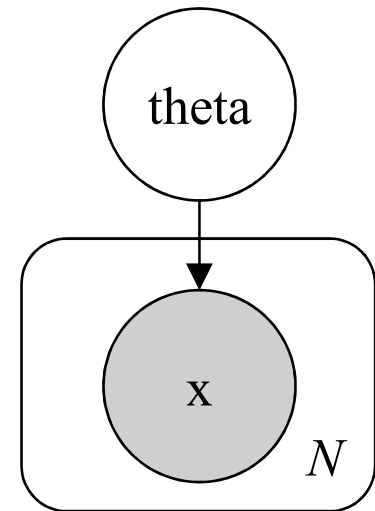
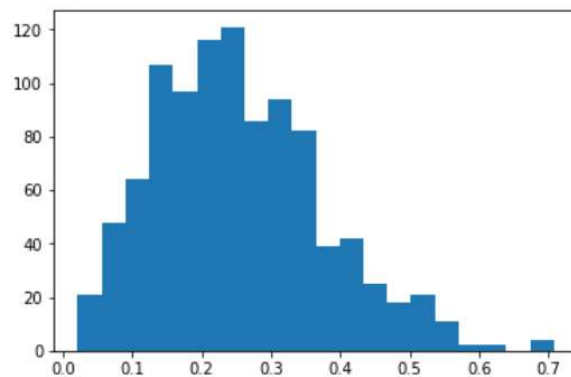
PyStan Example

```
In [2]: model = Model('''  
1      data {  
2          int<lower=0,upper=1> x[10];  
3      }  
4      parameters {  
5          real<lower=0,upper=1> theta;  
6      }  
7      model {  
8          theta ~ uniform(0,1);  
9          for (i in 1:10)  
10             x[i] ~ bernoulli(theta);  
11      }  
12  ''')
```

```
In [3]: flips = np.array([0, 1, 0, 0, 0, 0, 0, 0, 0, 1])  
1      run = model.sample(data={'x': flips}, random_='seed=42')  
2
```

```
In [4]: theta = run['theta']  
1      plt.hist(theta, bins=20)  
2      print('mean of theta: {:.3f}'.format(theta.mean()))  
3
```

mean of theta: 0.254



Outline

Strings Considered
Harmful

Name Reflection:
Keras.na

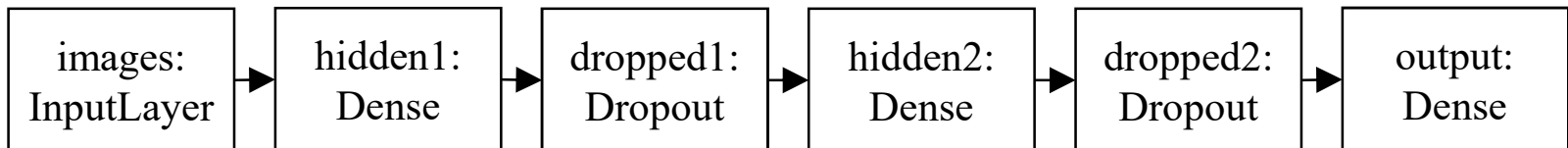
Reinterpreted
Python: YAPS

Keras with Fewer Strings

```
In [2]: import keras                                1
        from keras.models import Model              2
        from keras.layers import Dense, Dropout, Input 3
        keras_utils.wrap_layers(globals())          4
                                                5
        images = Input(shape=(N_PIXELS,))          6
        hidden1 = Dense(512, activation='relu')(images) 7
        dropped1 = Dropout(0.2)(hidden1)            8
        hidden2 = Dense(512, activation='relu')(dropped1) 9
        dropped2 = Dropout(0.2)(hidden2)           10
        output = Dense(N_LABELS, activation='softmax')(dropped2) 11
                                                12
        model = Model(inputs=images, outputs=output) 13
```

```
In [3]: keras_utils.plot_model(model)
```

Out[3]:



```
In [5]: model.compile(loss='categorical_crossentropy', 1
        optimizer='RMSprop',                          2
        metrics=['accuracy'])                          3
        model.fit({images: x_train}, {output: y_train}, 4
        batch_size=128, epochs=3,                     5
        validation_data=(x_test, y_test))              6
```

Name Reflection

```
1 import traceback
2 tb = traceback.extract_stack()
3 file_name, line_number, function_name, text = tb[-2]
4 import ast
5 tree = ast.parse(text, file_name)
6 assert isinstance(tree, ast.Module)
7 if isinstance(tree.body[0], ast.Assign):
8     targets = tree.body[0].targets
9     if len(targets) == 1 and isinstance(targets[0], ast.Name):
10         name = targets[0].id
```

Naming Heuristics

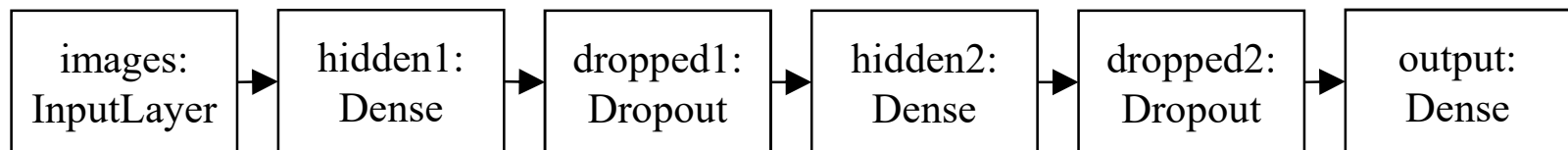
Case	Example	Name mapping
If explicit name argument	<code>aux_out = Dense(..., name='output2')(input)</code>	output2
Else if assigned to variable with unique name	<code>aux_out = Dense(...)(input)</code>	aux_out
Else	<code>h = Dense(...)(input)</code> <code>h = Dense(...)(h)</code>	h dense_2

Reverse Name Mapping

```
In [2]: import keras 1
        from keras.models import Model 2
        from keras.layers import Dense, Dropout, Input 3
        keras_utils.wrap_layers(globals()) 4
        5
        images = Input(shape=(N_PIXELS,)) 6
        hidden1 = Dense(512, activation='relu')(images) 7
        dropped1 = Dropout(0.2)(hidden1) 8
        hidden2 = Dense(512, activation='relu')(dropped1) 9
        dropped2 = Dropout(0.2)(hidden2) 10
        output = Dense(N_LABELS, activation='softmax')(dropped2) 11
        12
        model = Model(inputs=images, outputs=output) 13
```

```
In [3]: keras_utils.plot_model(model)
```

Out[3]:



```
In [5]: model.compile(loss='categorical_crossentropy', 1
        optimizer='RMSprop', 2
        metrics=['accuracy']) 3
        model.fit({images: x_train}, {output: y_train}, 4
        batch_size=128, epochs=3, 5
        validation_data=(x_test, y_test)) 6
```

- Names images, output are in scope (dict key does not start scope)
- Model wrapper re-keys the dictionary with strings from reflection

Keras.na Error Messages

```
In [1]: import keras                                1
        from keras.models import Model              2
        from keras.layers import Input, Conv2D      3
        import keras_utils                          4
        keras_utils.wrap_layers(globals())          5
                                                6
        img = Input(shape=(28,28))                 7
        conv = Conv2D(filters=10, kernel_size=3)(img) 8
        model = keras.models.Model(inputs=img, outputs=conv) 9
        model.compile(optimizer='sgd', loss='mean_squared_error') 10
```

Using TensorFlow backend.

ValueError: Input 0 is incompatible with layer conv: expected ndim=4, found ndim=3

Outline

Strings Considered
Harmful

Name Reflection:
Keras.na

Reinterpreted
Python: YAPS

Why Stan Probabilistic PL

- Widely adopted in the sciences
- Hundreds of examples
- Text books
- Dedicated conference
- High-level and self-contained language

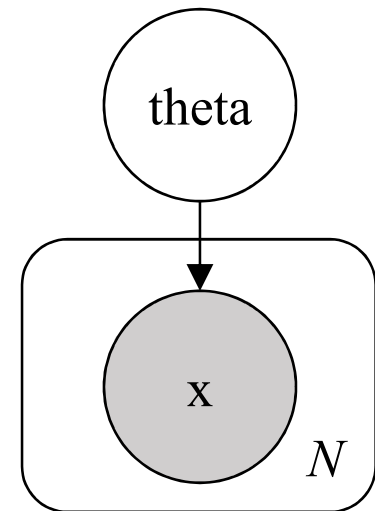
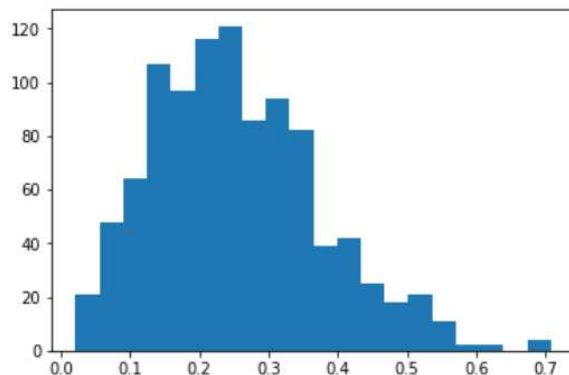
YAPS Example

```
In [2]: @yaps.model                                     1
        def coin(x: int(lower=0, upper=1)[10]):         2
            theta: real(lower=0, upper=1) <- uniform(0, 1) 3
            for i in range(1,11):                         4
                x[i] <- bernoulli(theta)                  5
```

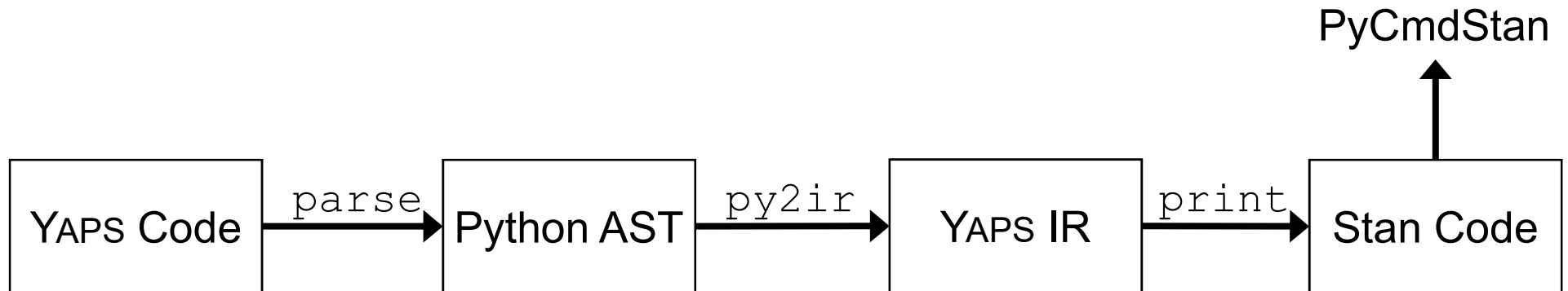
```
In [3]: flips = np.array([0, 1, 0, 0, 0, 0, 0, 0, 0, 1]) 1
        coin_obs = coin(x=flips)                          2
        coin_obs.sample(data=coin_obs.data, random_='seed=42') 3
```

```
In [4]: theta = coin_obs.posterior.theta                  1
        plt.hist(theta, bins=20)                          2
        print('mean of theta: {:.3f}'.format(theta.mean())) 3
```

mean of theta: 0.254



Reinterpreted Python



- User interacts with Yaps just like with any other Python library.
- Decorator `@yaps.model` replaces function, internally triggers Stan code generation.
- PyCmdStan takes care of the rest, invoking external Stan and C++ compilers.

Standard Python Tooling

```
@yaps.model
def coin(x: int(lower=0, upper=1),
        theta: real(lower=0, upper=1)) -> float:
    for i in range(1, 11):
        x[i] is bernoulli(thetap)
```

Undefined variable 'thetap'

[Quick Fix...](#) [Peek Problem](#)

Reverse Source Mapping

```
In [2]: @yaps.model                                     1
        def coin(x: int(lower=0, upper=1)[10]):         2
            for i in range(1,11):                       3
                x[i] <~ bernoulli(theta)                 4
```

```
In [3]: import numpy as np                             1
        flips = np.array([0, 1, 0, 0, 0, 0, 0, 0, 0, 1]) 2
        coin_obs = coin(x=flips)                        3
        coin_obs.sample(data=coin_obs.data)             4
```

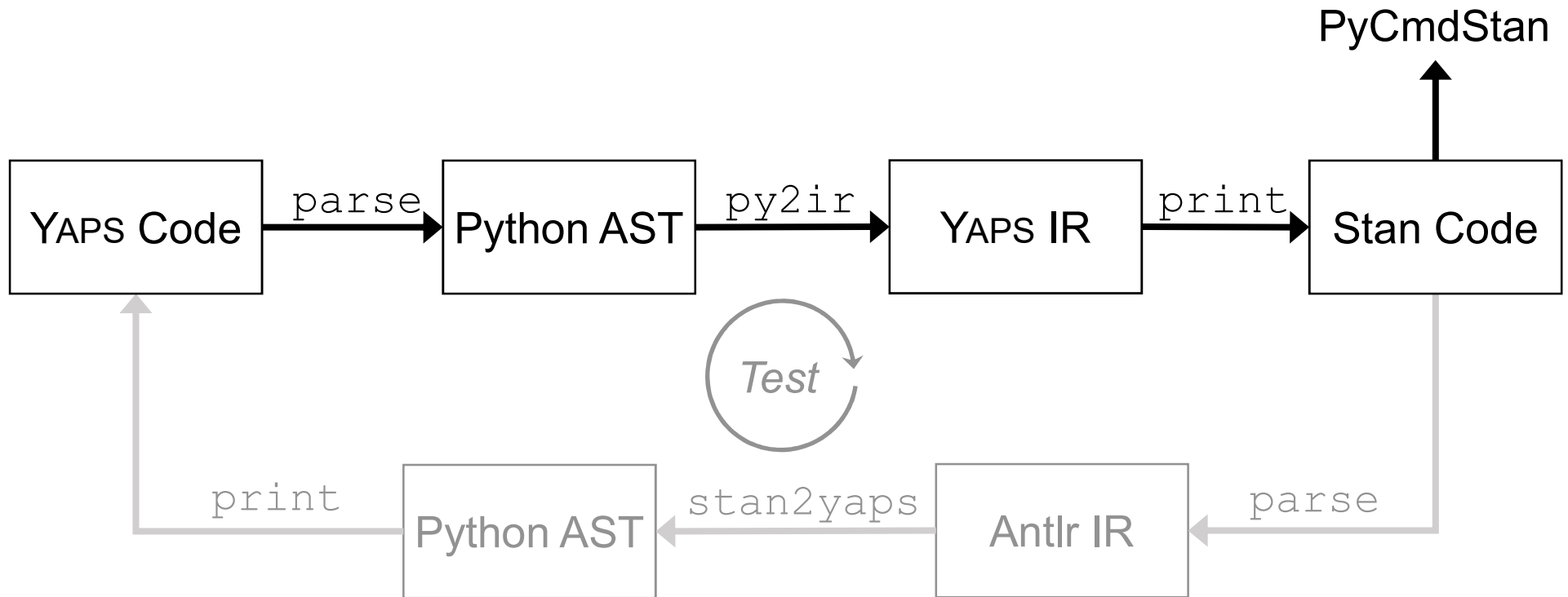
SYNTAX ERROR, MESSAGE(S) FROM PARSER:

variable "theta" does not exist.

error in '/Users/████████/.cache/pycmdstan/model-3aed19a9.stan'
' at line 4, column 27

```
-----
1: @yaps.model
2: def coin(x: int(lower=0, upper=1)[10]):
3:     for i in range(1,11):
4:         x[i] <~ bernoulli(theta)
                        ^
-----
```

Roundtrip Validation



	Total	Success
User manual	61	100%
Dev repo	721	97%
Examples repo	500	82%

The failed cases used deprecated Stan syntax.

YAPS Discussion

- Runtime validation
 - 13 example Stan models with datasets
 - Translated to YAPS, then back to Stan
 - Ran inference, found identical results
- Limitations
 - Typevar declaration for dependent types
 - Name clash with Python and Stan keywords
 - Debugger integration
- Open-source
 - <https://github.com/ibm/yaps/>
 - <https://pypi.org/project/yaps/>

Conclusions

Strings attached	No strings attached	Where to find
Keras	Keras.na	This paper
PyStan	YAPS	This paper; github
Sklearn	LALE	arXiv:1906.03957
Spark SQL	?	
?		