# If Not Notebooks, Then What?

Thoughts on notebooks and reproducibility

AAAI 2019 Workshop on Reproducible AI

Joel Grus @joelgrus



ALLEN INSTITUTE for ARTIFICIAL INTELLIGENCE

# About Me

- Research engineer on the <u>AllenNLP</u> team
- Previously SWE@Google, Data Science@VoloMetrix, ...
- Author of *Data Science from Scratch*
- Co-host of "<u>Adversarial Learning</u>" podcast
- The "Fizz Buzz in Tensorflow" guy
- Started "the first notebook war"

Allen







# Home About Programs Projects Careers Research Press

7 edd distant

# Al for the Common Good.

Our mission is to contribute to humanity through high-impact AI research and engineering.

# What is a Research Engineer?

- I'm not a researcher
  - (although I can pretend to be one)
- I'm an engineer who *cares deeply* about research
- My job is to help researchers be successful
  - by building tools for them
  - by advocating best practices
  - $\circ$  and also by collaborating with them
- And so my job involves thinking a lot about reproducibility



# Outline

- Why reproducibility?
- Why not notebooks?
- If not notebooks, then what?
- Reproducibility and AllenNLP
- Reproducibility and Beaker

# Why Reproducibility?

I take a somewhat expansive view of the role of reproducibility in science

# Reproducibility Helps With Correctness

If no one ever runs your code but you, it might be wrong



# Reproducibility Protects Against Bad Actors

- Certainly *you* would never try to publish fraudulent research, but what about your bitter rival?
- (Hopefully this is a minor consideration.)



# Reproducibility Makes It Easy to Try New Datasets

- Maybe your model works great on a different dataset
- Maybe it works terribly on a different dataset
- Hard datasets help move AI forward



## Reproducibility Makes It Easy to Try New Tasks

You haven't thought of (or tried) every way your model could be used

# Can you use BERT to generate text?

Just quickly wondering if you can use BERT to generate text. I'm using huggingface's pytorch pretrained BERT model (thanks!). I know BERT isn't designed to generate text, just wondering if it's possible. It's trained to predict a masked word, so maybe if I make a partial sentence, and add a fake mask to the end, it will predict the next word. As a first pass on this, I'll give it a sentence that has a dead giveaway last token, and see what happens.

### Reproducibility Enables Strong Baselines

Wouldn't you like your model to be the standard by which new models are judged?



# Reproducibility Is Necessary For Extensibility

"you can't stand on the shoulders of giants if they keep their shoulders private"



# Extensibility Leads to Progress



## Extensibility Leads to Progress



# **Fundamental Premise:**

The tools you choose and the processes you adopt can make reproducibility either a lot harder or a lot easier.



# What are Jupyter Notebooks?

- Computational environment that uses the lab notebook as a metaphor
- Easy to mix marked-up text, executable code, results, visualizations, interactive widgets, and so on
- Very popular with data scientists, less so (?) with researchers
- Often promoted as beneficial for "reproducibility"





#### Joel Grus @joelgrus

# slides for my "I Don't Like Notebooks" #JupyterCon talk:



I Don't Like Notebooks - Joel Grus - #JupyterCon ...

I Don't Like Notebooks hi, I'm Joel, and I don't like notebooks Joel Grus (@joelgrus) #JupyterCon 2018

docs.google.com

2.0K

9:55 AM - 24 Aug 2018 from Manhattan, NY

677

V

677 Retweets 1,985 Likes

136



111



# How Could Someone Not Like Notebooks?

- I had a lot of complaints (but also some positives and suggestions)
- It's a very comprehensive talk, you should check it out
- A lot of the complaints were around user-unfriendliness + confusingness
- One dislike was particularly relevant to this workshop:



NOTEBOOKS HINDER REPRODUCIBLE + EXTENSIBLE SCIENCE



# Notebooks for Reproducibility?



# Code That's in Notebooks Is Hard to Re-Run

### Notebooks Allow Out-of-Order Execution

CJUPYTEr madness			
File Edit View Insert Cell Kernel Help	Python 3 O		
B + ≫ 4  A < ↓ > A  Code	CellToolbar		
<pre>In [1]: def f(x): return x + 2</pre>			
In [3]: y = f(2)			
In [4]: y == 4			
Out[4]: False			
<pre>In [5]: print(y)</pre>			
5			

unless you are very disciplined, *you* may not even be able to reproduce your work

#### Notebooks Hard-Code Their Parameters

```
logits flatten = binary logits.view(-1, 2)
                return logits flatten, hidden
In [4]: rnn = RNN(input size=88, hidden size=512, num classes=88)
        rnn = rnn.cuda()
        criterion = nn.CrossEntropyLoss().cuda()
        criterion val = nn.CrossEntropyLoss(size average=False).cuda()
        learning rate = 0.001
        optimizer = torch.optim.Adam(rnn.parameters(), lr=learning_rate)
In [6]: %matplotlib notebook
        import sys, os
        sys.path.append("/home/
                                     /repos/pytorch-segmentation-detection/")
                                        .7repos/pytorch-segmentation-detection/vision/')
        sys.path.insert(0, '/nome/
```

# Notebooks Don't Help You Enforce Dependencies

Branch: master -			Find file	Сору	path
	/ notebooks /	n.ipy	nb		
added the link to dow	vnload midi library		11417b	5 on Ja	n 27
1 contributor					
1246 lines (1245 sloc) 104 KB		Raw Blame H	istory	<b>d</b> an t	Î
In [1]: import pretty_midi import numpy as np import torch import torch.nn as nn from torch.autograd imp import torch.utils.data import os import random	<b>port</b> Variable <b>a as data</b>				

# Notebooks Complicate Collaboration



# Collaboration is a Forcing Function for Reproducibility

### Notebooks Don't Play Well With Source Control

distributed source control is the de facto way to collaborate on coding projects

	1,	105	1,
	"source": [	106	"source": [
	"rng = np.random.RandomState(0)\n",	107	"rng = np.random.RandomState(0)\n",
114	— "for marker in ['o', ',', ',', 'x', '+', 'v', 'o', '8', '8', '8', '0']:\n",	108	+ "for marker in ['o', '≩', '+', 'v', '≤', '≥', '%', '₫', '≤']:\n",
115	<pre>" plt.plot(rng.rand(5), rng.rand(5), marker,\n",</pre>	109	<pre>" plt.plot(rng.rand(5), rng.rand(5), marker,\n",</pre>
116	<pre>- " label=\"marker='{0}'\".format(marker))\n",</pre>	110	<pre>+ " label=\"pointer='{0}'\".format(pointer))\n",</pre>
	<ul> <li>"plt.legend(numpoints=1)\n",</li> </ul>		<pre>+ "plt.legend(numpoints=3)\n",</pre>
118	<pre>- "plt.xlim(0, 1.8);"</pre>	112	+ "plt.xlim(0, 1.6);"
119	1	113	1
120	},	114	},
	{	115	(
牵	@@ -127,16 +121,14 @@		
	},	121	},
128	{	122	{
129	"cell_type": "code",	123	"cell_type": "code",
130	- "execution_count": 4,	124	+ "execution_count": 7,
	- "metadata": {	125	+ "metadata": {},
132	- "collapsed": false		
133	- },		
134	"outputs": [	126	"outputs": [
135	{	127	{
136	"data": {	128	"data": {
137	- "image/png":	129	+ "image/png":
	"iVBORw0KGgoAAAANSUhEUgAAAXoAAAD/CAYAAAD/qh1PAAAABHNCSVQICAgIfAhkiAAAAAlwSFl		"iVBORw0KGgoAAAANSUhEUgAAAXkAAAD0CAYAAAB+WlaPAAAABHNCSVQICAgIfAhkiAAAAAl
	z\nAAALEgAACxIB0t1+/AAAIABJREFUeJzt3XtcV0W6B/DfcBUZFE28VjjifWumeEMCGXUUuQgkK		zAAALEgAACxIB0t1+/AAAADl0RVh0U29mdHdhcmUAbWF0cGxvdGxpYiB2ZXJzaW9uIDIuMS4
	jgD\nlp6y2qmZHT1ak+nk9tLJTrX1aLl3KqCIN25K0gaiucUL3sp7Mt7AwhQVUBCYdf7owHbiPrd		odHRw0i8vbWF0cGxvdGxpYi5vcmcvNQv5yAAAIABJREFUeJzt3Xl4FFX28PFvJwSSsCmKoiy
	3rTXP\n9/PpD9Ywa/1YDY+LZ73vuyQcx3EghBAiWnasAxBCCLEsKvSEECJyV0gJIUTkqNATQojIU		AikNEQBBZBFkiWVC2iCxi/Dk6Kj0KjGFEfQi4oIi0jkZFwySBsIWwBRURRBSBj00ghKuoCfu
	aEnhBCR\no0JPCCEiZ1KhP3v2LGJiYmptz8zMRGRkJKKiorB9+3ZTDkEIIcREDsa+cc0GDUhJSYG		C0vX+kYQ3ZCHQW1VXn8/z5DF9q6rrXLo9qT59616HYRgIIYSwpyCzAxBCC0E9kuSFEMLGJMk
	rq6vB9srK\nSqxYsQK7du2Cs7Mzoq0jMXr0aLRt29bksIQQQprP6Ct6T09PrFmzptb2a9euwdPTE		NSZIXQggbkyQvhBA2JkleCCFsrJbZAVSUl5cnYzqFEMIF0dHRjoptlkvyANHR0S4dl5+fT8u
	1KpFI60jvD2\n9saJEydMCkkIIcR4Rhd6hUIBe3v7WttLSkrg5uZW87WrqyuKi4uNPQwhhBATmf1		cjbVJnw0D9DkwuNPnvLy8KtulXC0EEDYmSV4IIWxMkrwQQtiYJHkhhLAxt754VUq1B17SWne
	mrFQqRUlJSc3X\npaWlaNWqlbkPQwqhpImM7tFX+/NS0V5eXrhx4wYePXqEFi1a4MSJE5gxY0ad7		38CxwCpiutX5PKRUGpA0XAYeB4Vrrfe6cXwghxLm5fCWvlBoLvA+EVmgPAaYCvYA7gCSlVBP
	83NzTX18IQQYp08\nvb2b/L@mF3qJRAIASE9Px5MnTzBp0iQsXLqQ06dPB8dxmDRpEtq3b2+WsCy		j1vp2YAYw3tVzCvHsJyMjq8jISIKCqoiMjGTx4sVmh2QL7pRrfqYSqmhvCWzRWh/QWp8AvqR
	pVCok1CTU2u7o6Tou1aYhrA2d9dbda9abu709zd2aWVapW70u1a2v+Yv5C7eV1ce9E8ATa27evP7		Duga 2.0V. dobe fee in the TOM: web ToT. MINKCok a BODMCost OT 775 - 1/070HuE web/a DUg

# Notebooks Lock You Into Using Notebooks

#### The Module Finder

The finder is a simple object that tells you whether a name can be imported, and returns the appropriate loader. All this one does is check, when you do:

#### import mynotebook

it checks whether whether whether exists. If a notebook is found, then it returns a NotebookLoader.

### **Register the hook**

Now we register the NotebookFinder with sys.meta\_path

[ ]: sys.meta\_path.append(NotebookFinder())

#### Notebook Loader

**Totebooks a**:

ant to import code fro

a not plain Python files

Here we have our Notebook Loader. It's actually quite simple - once we figure out the filename of the module, all it does is:

load the notebook document into memory
 create an empty Module

3. execute every cell in the Module namespace

Since IPython cells can have extended syntax, the IPython transform is applied to turn each of these cells into their pure-Python counterparts before executing them. If all of your notebook cells are pure-Python, this step is unnecessary.

[]: class NotebookLoader(object): """Module Loader for Jupyter Notebooks""" def \_\_init\_\_(self, path=None): self.shell = InteractiveShell.instance() self.path = path

def load\_module(self, fullname):
 """import a notebook as a module"""
 path = find\_notebook(fullname, self.path)

print ("importing Jupyter notebook from %s" % path)

# load the notebook object
with io.open(path, 'r', encoding='utf-8') as f:
 nb = read(f, 4)

### Notebooks Conflate "Library" Code and "Experiment" Code

#### In [3]: class RNN(nn.Module):

def init (self, input size, hidden size, num classes, n layers=2):

super(RNN, self).\_\_init\_\_()

self.input\_size = input\_size
self.hidden\_size = hidden\_size
self.num\_classes = num\_classes
self.n\_layers = n\_layers

self.notes\_encoder = nn.Linear(in\_features=input\_size, out\_features=hidden\_size)

self.lstm = nn.LSTM(hidden\_size, hidden\_size, n\_layers)

self.logits\_fc = nn.Linear(hidden\_size, num\_classes)

def forward(self, input sequences, input sequences lengths, hidden=None):

```
batch_size = input_sequences.shape[1]
```

notes\_encoded = self.notes\_encoder(input\_sequences)

```
# Here we run rnns only on non-padded regions of the batch
```

```
packed = torch.mn.utils.rnn.pack_padded_sequence(notes_encoded, input_sequences_lengths)
outputs, hidden = self.lstm(packed, hidden)
outputs, output_lengths = torch.nn.utils.rnn.pad_packed_sequence(outputs) # unpack (back
to padded)
```

```
logits = self.logits_fc(outputs)
```

logits = logits.transpose(0, 1).contiguous()

```
neg_logits = (1 - logits)
```

# Since the BCE loss doesn't support masking, we use the crossentropy binary\_logits = torch.stack((logits, neg\_logits), dim=3).contiguous()

logits flatten = binary logits.view(-1, 2)

return logits\_flatten, hidden

In [4]: rnn = RNN(input\_size=88, hidden\_size=512, num\_classes=88)
rnn = rnn.cuda()

criterion = nn.CrossEntropyLoss().cuda()

criterion val = nn.CrossEntropyLoss(size average=False).cuda()

learning\_rate = 0.001
optimizer = torch.optim.Adam(rnn.parameters(), lr=learning\_rate)

In [8]:	: clip = 1.0 epochs_number = 12000000 sample_history = [] best_val_loss = float("inf")		
	<pre>for epoch_number in xrange(epochs_number):</pre>		
	for	batch in trainset_loader:	
		<pre>post_processed_batch_tuple = post_process_sequence_batch(batch)</pre>	
	tuple	<pre>input_sequences_batch, output_sequences_batch, sequences_lengths = post_processed_batch_</pre>	
	a() )	<pre>output_sequences_batch_var = Variable( output_sequences_batch.contiguous().view(-1).cud</pre>	
		<pre>input_sequences_batch_var = Variable( input_sequences_batch.cuda() )</pre>	
		optimizer.zero_grad()	
		<pre>logits, _ = rnn(input_sequences_batch_var, sequences_lengths)</pre>	
		loss = criterion(logits, output_sequences_batch_var) loss_list.append( loss.data[0] ) loss.backward()	
		<pre>torch.nn.utils.clip_grad_norm(rnn.parameters(), clip)</pre>	
		optimizer.step()	
	cur val	<pre>rent_val_loss = validate() _list.append(current_val_loss)</pre>	
	<pre>if current_val_loss &lt; best_val_loss:</pre>		
		<pre>torch.save(rnn.state_dict(), 'music_rnn.pth') best_val_loss = current_val_loss</pre>	

Notebooks are a recipe for poorly factored code



12.26 AM - 15 Jul 2018



@fchollet

Following

# Notebooks Frustrate Software-Engineering Best Practices

# Software Engineering?

- You may not think of your Al experiments as software engineering
- But I do
- And you should!
- I sometimes give talks on "why data scientists should care about software engineering best practices"
- This is sort of my "why AI researchers should care about software engineering best practices" variant



# Notebooks Complicate Code Reviews

- Code reviews are an early bulwark against incorrect code (and hence incorrect science)
- Code reviews help you grow as a coder and as a scientist



# Notebooks Make it Hard to Unit Test

- Unit tests assure you that small pieces of your logic are correct
- Unit tests make it easy to iterate by running your model end-to-end on small datasets
- Unit tests make it safe to refactor your code
- Unit tests are small working examples
- Yes, you can put asserts in your notebook, but you want to be able to run your tests without running your experiment



#	load data	
#	maybe assert	something
	t	
#	train model	
#	maybe assert	something
	-	-
#	get results	
#	mavbe assert	something
ш.		50112118
#	save results	

Yeah, but what do unit tests for Al experiments even look like?
### Unit Tests for AI Experiments

```
def test_forward_pass_runs_correctly(self):
```

```
training_tensors = self.dataset.as_tensor_dict() tiny known dataset
tags = output_dict['tags'] check that output has the right fields
assert len(tags) == 2
                          check that output has the right shape
assert len(tags[0]) == 7
assert len(tags[1]) == 7
for example_tags in tags:
   for tag_id in example_tags:
       tag = self.model.vocab.get_token_from_index(tag_id, namespace="labels")
       assert tag in {'0', 'I-ORG', 'I-PER', 'I-LOC'}
```

check that output has reasonable values

# Notebooks Make Science Harder

#### Notebooks Are Hard to Parametrize



# Notebooks Lump Together Code and Data

- Often you'll run many experiments with the same code but different parameters
- For analysis it can be good to have your results mixed in
- For science you want to aggregate results across experiments
- For science you want your results to have a life of their own



## Notebook Code Can't Easily Be Built On



## In Short

- "Notebooks as a source of reproducibility" presupposes a *static* view of Al as a science:
  - "I did some science, now here is an artifact containing the code and data"
- I'm advocating for a *dynamic* view:
  - "I did some science, now you do some more science *on top of it*"
  - Which is the kind of reproducibility that moves AI forward

# If Not Notebooks, Then What?

### Code in Modules

. . .

```
# models/crf tagger.py
class CrfTagger(Model):
   ......
  The ``CrfTagger`` encodes a sequence of text with a ``Seq2SeqEncoder``,
  then uses a Conditional Random Field model to predict a tag for each token in the sequence.
  def init (self, vocab: Vocabulary,
                text field embedder: TextFieldEmbedder,
                encoder: Seq2SeqEncoder,
                label namespace: str = "labels",
                feedforward: Optional[FeedForward] = None,
                label encoding: Optional[str] = None,
                include start end transitions: bool = True,
                constrain crf decoding: bool = None,
                calculate span f1: bool = None,
                dropout: Optional[float] = None,
                verbose metrics: bool = False,
                initializer: InitializerApplicator = InitializerApplicator(),
                regularizer: Optional[RegularizerApplicator] = None) -> None:
       super(). init (vocab, regularizer)
```

### Unit Tests

- develop your model on a tiny test set
- iterate
- run the tests frequently



# The best time to find mistakes is before you run your experiments!



#### **Be Explicit About Your Dependencies**

#### ESSENTIAL LIBRARIES FOR MAIN FUNCTIONALITY ####

# This installs Pytorch for CUDA 8 only. If you are using a newer version, # please visit http://pytorch.org/ and install the relevant version. # For now AllenNLP works with both PyTorch 1.0 and 0.4.1. Expect that in # the future only >=1.0 will be supported. torch>=0.4.1

# Parameter parsing (but not on Windows).
jsonnet==0.10.0 ; sys.platform != 'win32'

# Adds an @overrides decorator for better documentation and error checking when using subclasses. overrides

# Used by some old code. We moved away from it because it's too slow, but some old code still # imports this. nltk

```
# Mainly used for the faster tokenizer.
spacy>=2.0,<2.1</pre>
```

### Make It Easy To Vary Parameters

```
export BERT_BASE_DIR=/path/to/bert/uncased_L-12_H-768_A-12
export GLUE_DIR=/path/to/glue
```

```
python run_classifier.py \
```

```
--task_name=MRPC \
```

```
--do_train=true \
```

```
--do_eval=true \
```

```
--data_dir=$GLUE_DIR/MRPC \
```

```
--vocab_file=$BERT_BASE_DIR/vocab.txt \
```

```
--bert_config_file=$BERT_BASE_DIR/bert_config.json \
```

```
--init_checkpoint=$BERT_BASE_DIR/bert_model.ckpt \
```

```
--max_seq_length=128 \
```

```
--train_batch_size=32 \setminus
```

```
--learning_rate=2e-5 \setminus
```

```
--num_train_epochs=3.0 \
```

```
--output_dir=/tmp/mrpc_output/
```

## Consider Docker Images

- Create container with OS + environment + code (+ data?)
- Can share and anyone can run (in theory)
- Build up step by step with smart caching when you change it



root@7d1f120a83e9:/stage/allennlp# echo '{"sentence": "Did Uriah honestly think he could beat the
game in under three hours?"}' | allennlp predict
https://s3-us-west-2.amazonaws.com/allennlp/models/srl-model-2018.05.25.tar.gz -

#### **Provide Instructions**

#### Installation

This repo was tested on Python 3.5+ and PyTorch 0.4.1/1.0.0

#### With pip

PyTorch pretrained bert can be installed by pip as follows:

pip install pytorch-pretrained-bert

#### From source

Clone the repository and run:

pip install [--editable] .

A series of tests is included in the tests folder and can be run using pytest (install pytest if needed: pip install pytest).

You can run the tests with the command:

python -m pytest -sv tests/

# Reproducibility and AllenNLP

#### What is AllenNLP?



#### Deep learning for NLP

AllenNLP makes it easy to design and evaluate new deep learning models for nearly any NLP problem, along with the infrastructure to easily run them in the cloud or on your laptop.

#### State of the art models

AllenNLP includes reference implementations of high quality models for both core NLP problems (e.g. semantic role labeling) and NLP applications (e.g. textual entailment).

#### **Get Started**

#### **View Models**

### Programming to Higher-Level Abstractions

. . .

```
# models/crf tagger.py
class CrfTagger(Model):
   ......
   The ``CrfTagger`` encodes a sequence of text with a ``Seq2SeqEncoder``,
   then uses a Conditional Random Field model to predict a tag for each token in the sequence.
   def __init__(self, vocab: Vocabulary,
                text field embedder: TextFieldEmbedder,
                encoder: Seq2SeqEncoder,
                label namespace: str = "labels",
                feedforward: Optional[FeedForward] = None,
                label encoding: Optional[str] = None,
                include start end transitions: bool = True,
                constrain crf decoding: bool = None,
                calculate span f1: bool = None,
                dropout: Optional[float] = None,
                verbose metrics: bool = False,
                initializer: InitializerApplicator = InitializerApplicator(),
                regularizer: Optional[RegularizerApplicator] = None) -> None:
       super(). init (vocab, regularizer)
```

### **Declarative Configuration**

```
"model": {
 "type": "crf_tagger",
  "label encoding": "BIOUL",
  "dropout": 0.5,
  "include start end transitions": false,
  "text field embedder": {
    "token embedders": {
      "tokens": {
          "type": "embedding",
          "embedding dim": 50,
          "pretrained file": "/path/to/glove.txt.gz",
          "trainable": true
     },
      "elmo":{
          "type": "elmo token embedder",
          "options file": "/path/to/elmo/options.json",
          "weight file": "/path/to/elmo/weights.hdf5",
          "do layer norm": false,
          "dropout": 0.0
     },
```

```
"token characters": {
        "type": "character_encoding",
        "embedding": {
            "embedding dim": 16
        },
        "encoder": {
            "type": "cnn",
            "embedding dim": 16,
            "num filters": 128,
            "ngram filter sizes": [3],
            "conv layer activation": "relu"
},
"encoder": {
    "type": "lstm",
    "input size": 1202,
    "hidden size": 200,
    "num layers": 2,
    "dropout": 0.5,
    "bidirectional": true
},
"regularizer": [
    "scalar parameters",
      "tvpe": "12",
      "alpha": 0.1
```

},



#### **Command-Line Tools**

#### Using the Command Line Tool

In fact, we provide a command line tool that handles most common AllenNLP tasks, so in practice you probably wouldn't read the params or call train\_model yourself (although you can), you'd just run the command

\$ allennlp train tutorials/tagger/experiment.jsonnet \
 -s /tmp/serialization\_dir \
 --include-package tutorials.tagger.config\_allennlp

#### Interactive Demos and Visualization of Model Internals

#### **Allen**NLP

#### Machine Comprehension

#### **Textual Entailment**

Semantic Role Labeling

Coreference Resolution

Named Entity Recognition

Constituency Parsing

Your model here!

#### Textual Entailment

Textual Entailment (TE) takes a pair of sentences and predicts whether the facts in the first necessarily imply the facts in the second one. The AllenNLP toolkit provides the following TE visualization, which can be run for any TE model you develop. This page demonstrates a reimplementation of the decomposable attention model (Parikh et al, 2017), which was state of the art for the SNLI benchmark (short sentences about visual scenes) in 2016. Rather than pretrained Glove vectors, this model uses ELMo embeddings, which are completely character based and improve performance by 2%

Enter text or Choose an example... 🗘

#### Premise

If you help the needy, God will reward you.

#### Hypothesis

Giving money to the poor has good consequences.

RUN



#### Model internals (beta)

#### premise to hypothesis attention

hypothesis to premise attention

For every premise word, the model computes an attention over the hypothesis words. This heatmap shows that attention, which is normalized for every row in the matrix.



## Higher-Level NLP Abstractions as Library Primitives

- Field + Instance (nice representation of examples)
  - o question\_field = TextField("What is the ...", token\_indexers)
  - o instance = Instance({"question": question\_field, "passage": passage\_field)
- Vocabulary (map word <-> index, label <-> index, etc)
- TokenIndexer (map token -> [index1, ...., indexn]
  - could be one per word
  - could be one per character
  - could be one per wordpiece
- TokenEmbedder (map indices -> embedding vectors)
- Seq2VecEncoder (map [v1, ..., vn] -> w)
- Seq2SeqEncoder (map [v1, .., vn] -> [w1, ..., wn])
- ...and many more

"I want to try using BERT vectors instead of GloVe vectors"



"Oh, great, now I'm going to make lots nges of my code aintain ferent all versions so that my results are reproducible"

"I'll just make a new config file for the BERT version!"

```
"token_indexers": {
    "tokens": {
        "type": "single_id",
        "lowercase_tokens": true
    },
```

```
"token_characters": {
    "type": "characters",
    "min_padding_length": 3
}
```

}

```
"token_indexers": {
    "bert": {
        "type": "bert-pretrained",
        "pretrained_model": std.extVar("BERT_VOCAB"),
        "do_lowercase": false,
        "use_starting_offsets": true
    },
    "token_characters": {
        "type": "characters",
        "min_padding_length": 3
    }
}
```

```
"text field embedder": {
  "token embedders": {
    "tokens": {
        "type": "embedding",
        "embedding_dim": 50,
        "pretrained_file": "/path/to/glove.tar.gz"
        "trainable": true
    },
    "token characters": {
        "type": "character_encoding",
        "embedding": {
            "embedding_dim": 16
        },
        "encoder": {
            "type": "cnn",
            "embedding_dim": 16,
            "num filters": 128,
            "ngram_filter_sizes": [3],
            "conv_layer_activation": "relu"
   },
                                                       },
},
```

```
"text field embedder": {
    "allow_unmatched_keys": true,
    "embedder to indexer map": {
        "bert": ["bert", "bert-offsets"],
        "token characters": ["token_characters"],
    },
    "token embedders": {
        "bert": {
            "type": "bert-pretrained",
            "pretrained_model": std.extVar("BERT_WEIGHTS")
        },
        "token_characters": {
            "type": "character_encoding",
            "embedding": {
                "embedding_dim": 16
            },
            "encoder": {
                "type": "cnn",
                "embedding_dim": 16,
                "num_filters": 128,
                "ngram_filter_sizes": [3],
                "conv_layer_activation": "relu"
        }
```

## "Want to Reproduce My Results?"

- create a virtual environment
- clone my GitHub repo and pip install its dependencies
  - including a specific version of AllenNLP
  - which includes a specific version of PyTorch
- each experiment has its own JSON configuration file
- allennlp train specific\_experiment.json -s /tmp/results

# Reproducibility and Beaker

## What is Beaker?

- Kubernetes-based platform for rapid experimentation
- Specify experiments as Docker containers + config.yaml
- Request GPUs + Memory + etc
- Upload or mount existing datasets
- Track results
- (Currently runs on GKE, working on a "runs on your machines" version)
- Mostly internal-only for now (sorry)





### What is Beaker?

beaker experiment run  $\$ 

- --name wordcount-moby  $\backslash$
- --blueprint examples/wordcount \
- --source examples/moby:/input \
- --result-path /output

ex_ndcfzinlot	02 🖉			C <sup>4</sup> Re-run	Program Arguments	python -m allen	nlp.run train /config.ison	-s /outputfile-frie	ndly-logging
NER with char + BERT (using	[SEP] and [CLS])	descripti	ion		Environment Variables				,
Created 2 months ago by joel	Ig				Environment variables	REPT VOCAR		value	
Tasks Experiment Gra	aph Spec Permission	IS				DERT_VOCAD		/data/bert-vocab	
						BERI_WEIGHTS		/data/bert-weights	
C training						NER_TEST_A_PATH		/conll2003/eng.testa	
Uddining	training					NER_TEST_B_PATH		/conll2003/eng.testb	
	Succeeded Finished 2 months ago. Ran for 2 hours. Created by: joelg				narameters	NER_TRAIN_DATA_PATH		/conll2003/eng.train	
					parameters	dataset_reader_coding_scheme		BIOUL	
	Description	None 🖉 logs				dataset_reader_tag_label		ner	
					dataset_reader_token_indexers_bert_do_lowerc dataset_reader_token_indexers_bert_pretrained dataset_reader_token_indexers_bert_type		False /data/bert-vocab bert-pretrained		
	Logs	2506 Lines (377.3KiB) Show Full Log Download Full Log							
	Results	ds_og42u9uwhhld							
					Show all				
	Metrics	{     "best_epoch": 37,     "best_validation_accuracy": 0.9904209337642615,     "best_validation_accuracy": 0.90152020327642615,		Requirements	1 GPU 0.0B Memory				
		"best_validation_fl.messure-overall": 0.95274928535391, "best_validation_fl.messure-overall": 0.95274928535391, "best_validation_loss": 51.06252528171913,			Blueprint	bp_w3b6mh1llks8 docker image			
		"best_validation_precision-overall": 0.9519489247311828, "best_validation_recall-overall": 0.9535509929316729, "epoch": 61,		Experiment	ex_ndcfzinlot02				
metrics		"training_accuracy": 0.998340053334381, "training_accuracy3": 0.9985414078115715, "training_duration": "01:52:06", "training_epochs": 61,		Input Datasets	Dataset ID bert-base-cased-voca	Mount Path b /data/bert-vocab			
		"training_fl-m "training loss	easure-overall": 0.990996360230 ": 3.9061280575665562,	6816,		bert-base-cased	/data/bert-weights	datasets	;
		"training_prec	ision-overall": 0.9913550804871	817,		conll2003	/conll2003		
		"training_teca "training_star "validation_ac	t_epoch": 0, curacy": 0.9893111638954869,	2		ds_5oc1fv432odm	/config.json		
		"validation_accuracy3": 0.9906740391729294, "validation_f1-measure-overall": 0.9481381860972855, "validation_loss": 65.0678040747549,			ID	tk_63w97m1997ba			
		"validation_precision-overall": 0.947103274559194, "validation_recall-overall": 0.9491753618310333 }		9194, 33	Cost	\$4.887	cost		

### Organize Experiments Into Groups

ideally you'd give them more descriptive names, though

Compa	risons							
Q Sea	arch by task or experir	nent name				Search	111 Manage Comparisons	🛓 Export to CSV
•	Task	♦ Ex	periment	Ŧ	best_epoch	\$	best_validation_accuracy	\$
S	tk_8iw1xf6p8jrs	ex	_znnqknv5o3jm		4.0		0.764	
0	tk_xyf3kctrkg0j	ex	_z58l8spierr1		4.0		0.767	
S	tk_1dymlce34chv	ex	yk6k1m5a7i6i		.3.0		0.766	
Ø	tk_6d2p2kfr2ysf	ex	y0lzev942p28		4.0		0.75	
S	tk_hzwsz4issg72	ex	_xkrg635uyrg4		4.0		0.767	
S	tk_vovo4vp329ya	ex	_xesvw1l8s3um		9.0		0.762	
S	tk_zofxyhjc7p14	ex	_x8sq9ax7cqk6		4.0		0.765	
S	tk_whsyiaj04zab	ex	_x6eczhvm35d7		4.0		0.767	
S	tk_xvirkb29ky4j	ex	_wqspqyc4138c		4.0		0.763	
S	tk_7schpvcy3m6z	ex	_wq2feuiikmna		5.0		0.741	
S	tk_9wlpjr3k6vbo	ex	_w6ngi069j1gn		9.0		0.741	
S	tk_78lm112dptkh	ex	_w2qmp3ljj4vt		4.0		0.753	

## Beaker and Reproducibility

- old code + new data => upload the dataset, reuse the blueprint
- new code + old data => create the blueprint, point at existing dataset
- want to see previous results?
  - inputs + logs + outputs stored "forever"
  - record of every experiment run + results
  - share with a link

# To Sum Up

- Reproducibility is important for more than the obvious reasons
- Your choices of tools and processes make reproducibility easier or harder
- There are several ways that notebooks make reproducibility harder
- Search out tools that make reproducibility easier
- Adopt processes that make reproducibility easier

### A Few Related Presentations

• I Don't Like Notebooks

The talk that launched a thousand arguments. Heavy on the memes.

Writing Code for NLP Research

EMNLP 2018 tutorial from me + Matt Gardner + Mark Neumann, goes much deeper into "what good research code looks like"

How Becoming Not a Data Scientist Made Me a Better Data Scientist

Explores some similar themes in the context of "why data scientists should care about software engineering best practices"
## Any questions?

me: <u>@joelgrus</u>

Al2: allenai.org

AllenNLP: allennlp.org

Beaker\*: beaker.org



will tweet out slides from

@joelgrus and @ai2\_allennlp