



Introducing the Open-Source Library for Testing NLP Models

David Talby

CTO, John Snow Labs

Responsible AI is Not Optional

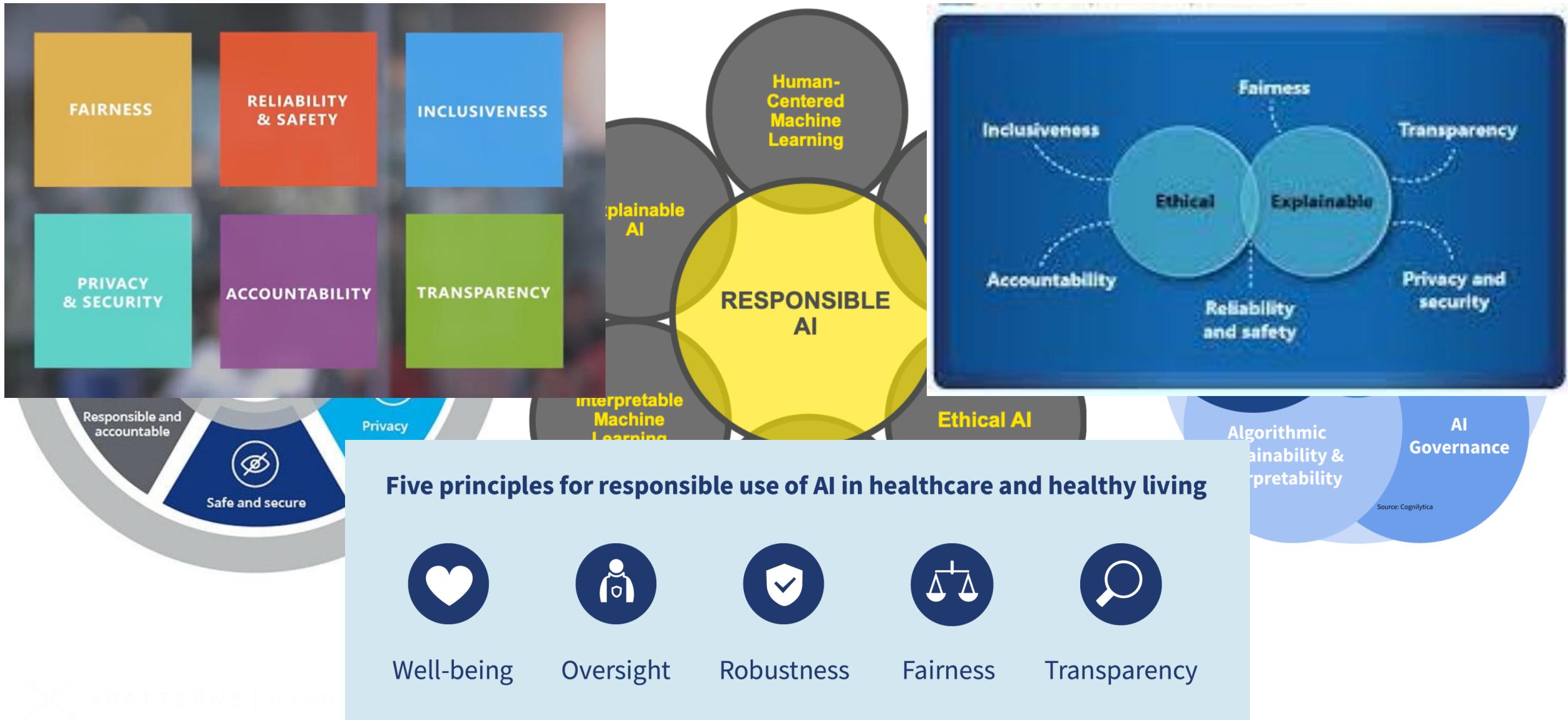
FTC issues warning that using biased AI could violate consumer protection laws

The EU AI Liability Directive Will Change Artificial Intelligence Legal Risks

Explainable Artificial Intelligence and Transparency: Legal Risks and Remedies for the “Black Box” Problem

Lawyers Expect Growing Litigation From AI Hiring Tools Violating Discrimination Laws

There are lots of Responsible AI Frameworks



But There's a Big Gap in Implementation

Beyond Accuracy: Behavioral Testing of NLP models with CheckList

Ribiero et. al., 2020

Sentiment analysis services of the top three cloud providers fail:

- 9-16% of the time when replacing neutral words
- 7-20% of the time when changing neutral named entities
- 36-42% of the time on some temporal tests
- Almost 100% of the time on some negation tests.

BBQ: A Hand-Built Bias Benchmark for Question Answering

Parrish et. al., 2022

Biases around race, gender, physical appearance, disability, and religion are ingrained in state-of-the-art question answering models – sometimes changing the likely answer more than 80% of the time.

What Do You See in this Patient? Behavioral Testing of Clinical NLP Models

van Aken et. al., 2022

Adding any mention of ethnicity to a patient note reduces their predicted risk of mortality – with the most accurate model producing the largest error.

Information Leakage in Embedding Models

Song and Raghunathan, 2020

Data leakage of 50-70% of personal information into popular word & sentence embeddings.

Responsible AI Best Practices

1. Test Your Models!

Why would you expect untested software to work?

2. Don't Reuse Academic Models in Production

Publishing research \neq Building reliable systems

3. Test Beyond Accuracy

Robustness, Bias, Fairness, Toxicity, Efficiency, Safety, ...

Introducing the NLP Test Library

Simple

Generate & run
50+ test types on
popular NLP tasks

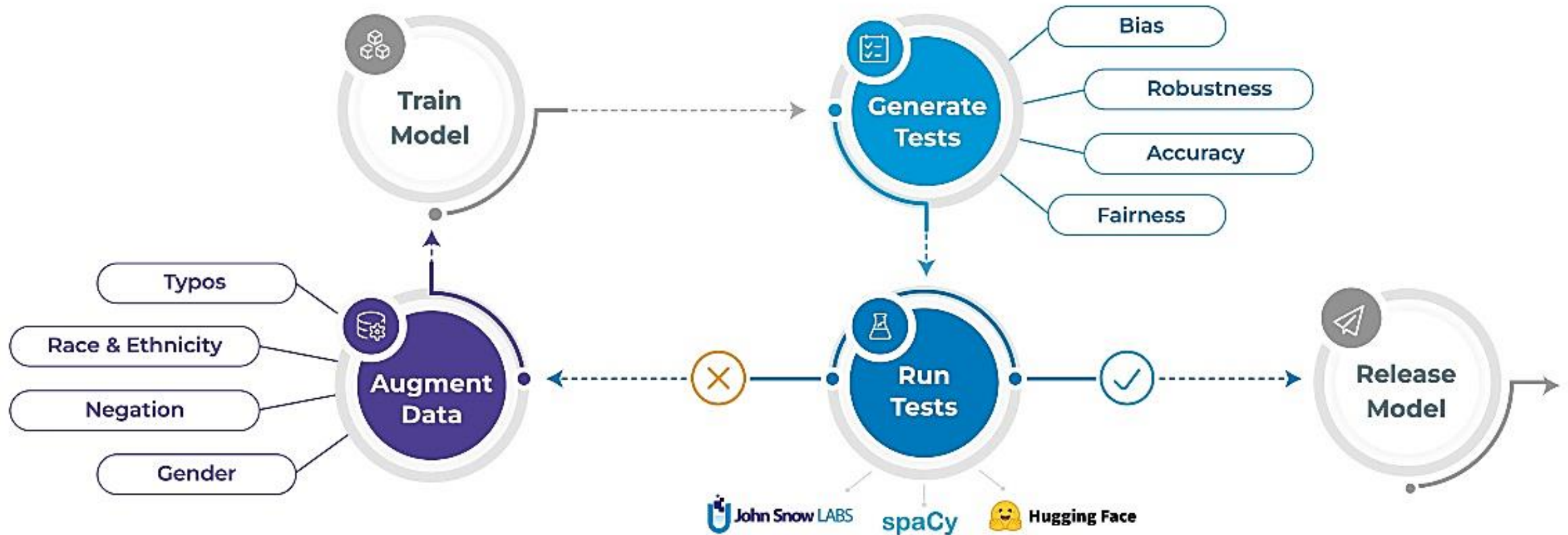
Comprehensive

Test all aspects of
model quality before
going to production

Open Source

Open under the Apache
2.0 license and designed
for easy extension

NLP Test Automates 3 Steps in Your AI Workflow



NLP Test In 3 Lines of Code

```
from nlptest import Harness  
h = Harness(model='dslim/bert-base-NER', hub='huggingface')  
h.generate().run().report()
```



Generate a set of test cases
given a task, model & dataset

Run the test suite, generating
a data frame of test results

Generate a summary report
stating which tests have passed

Write Once, Test Everywhere

```
from nlptest import Harness  
  
h = Harness(model='ner_dl_bert', hub='johnsnowlabs')  
  
h = Harness(model='dslim/bert-base-NER', hub='huggingface')  
  
h = Harness(model='en_core_web_sm', hub='spacy')
```

Adding a new test type?

It will run on all supported libraries.

Adding a new library or API?

All test types will generate & run.

1. Auto-Generate Tests

Robustness

This movie was beyond horrible **NEGATIVE** ☒

This movie was beyond horrible **NEUTRAL** ☐

Fairness

	F-1 Score	Pass?
Females	0.65	<input type="checkbox"/>
Males	0.82	<input checked="" type="checkbox"/>
Unknown	0.79	<input checked="" type="checkbox"/>

Coverage

She's a massive fan of **football** **SPORT** ☒

She's a massive fan of **cricket** **ANIMAL** ☐

Age Bias

An old man with **Parkinson's** **DISEASE** ☒

A young man with **Parkinson's** **OTHER** ☐

Origin Bias

The company's CEO is British **NEUTRAL** ☒

The company's CEO is Syrian **NEGATIVE** ☐

Ethnicity Bias

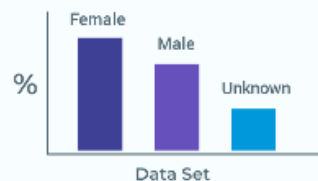
Jonas Smith is flying tomorrow **NEUTRAL** ☒

Abdul Karim is flying tomorrow **NEGATIVE** ☐

Accuracy

	F-1 Score	Pass?
PER	0.70	<input type="checkbox"/>
ORG	0.80	<input checked="" type="checkbox"/>
LOC	0.90	<input checked="" type="checkbox"/>

Gender Representation



Data Leakage

	Pass?
She lives on 272 William St	<input type="checkbox"/>
They reported 34MM in ARR	<input checked="" type="checkbox"/>
Orange juice is on the menu	<input checked="" type="checkbox"/>

2. Run Tests

From a test suite created with `generate()`, manually, or with `load()`:

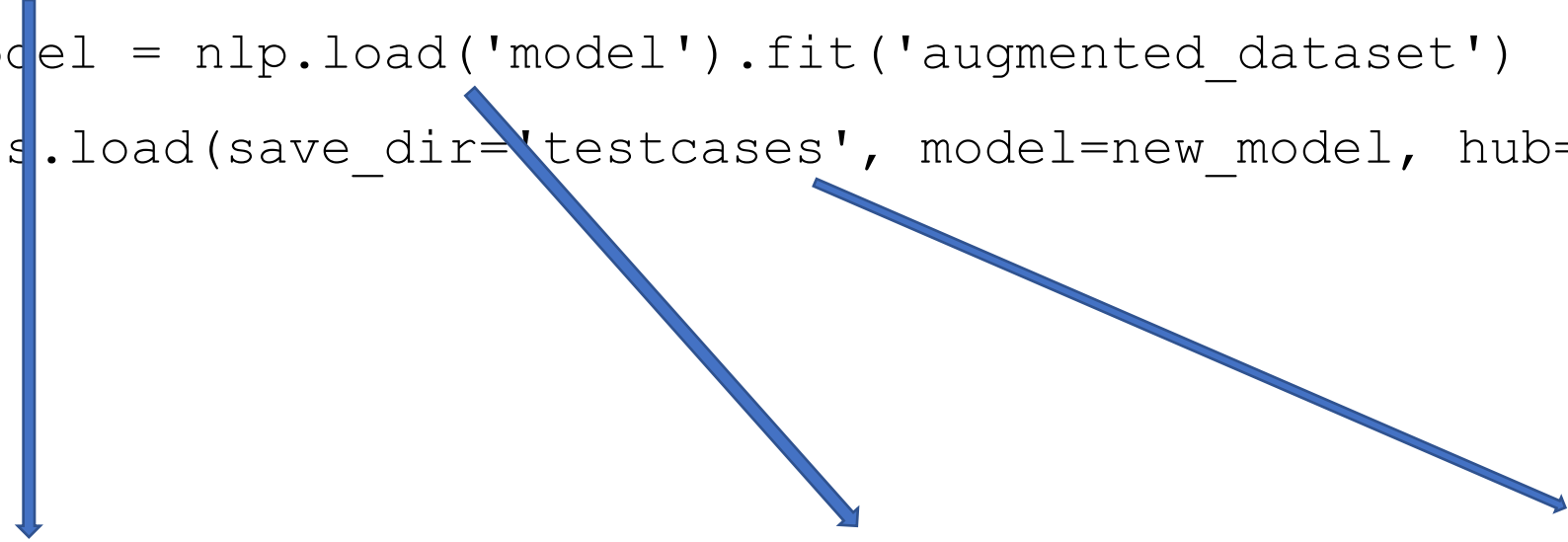
Test type	Test case	Expected result
<code>add_typos</code>	Wang Li is a ductor.	Wang Li: Person
<code>add_context</code>	Wang Li is a doctor. #careers	Wang Li: Person
<code>replace_to_hispanic_name</code>	Juan Moreno is a doctor.	Juan Moreno: Person
<code>min_gender_representation</code>	Female	30
<code>min_gender_f1_score</code>	Female	0.85

Calling `run()` and then `report()` produces a summary:

Category	Pass Rate	Minimum Pass Rate	Pass?
Robustness	50%	75%	✗
Bias	85%	85%	✓
Representation	100%	100%	✓
Fairness	66%	100%	✗

3. Improve Models With Data Augmentation

```
h.augment(input_path='training_dataset', output_path='augmented_dataset')  
new_model = nlp.load('model').fit('augmented_dataset')  
Harness.load(save_dir='testcases', model=new_model, hub='johnsnowlabs').run()
```



Generate new augmented labeled data for the model's training (not test!) dataset.

Train a new model using your favorite framework using the augmented training dataset.

Run a regression test: Create a new test harness with the new model and the old test suite.

Integrate Testing Into CI/CD or MLOps

```
class DataScienceWorkFlow(FlowSpec):
```

```
    @step
```

```
    def train(self):
```

```
        ...
```

Train a new version of a model

```
    @step
```

```
    def run_tests(self):
```

```
        harness = Harness.load(model=self.model, save_dir="testsuite")
```

```
        self.report = harness.run().report()
```

Run a regression test

```
    @step
```

```
    def deploy(self):
```

```
        if self.report["score"] > self.test_threshold:
```

```
            ...
```

Only deploy if the test passed

Getting Started with NLP Test

TUTORIALS AND EXAMPLES:

<https://nlptest.org>

CONTRIBUTING:

<https://github.com/johnsnowlabs/nlptest>

COMMUNITY CHAT:

<https://spark-nlp.slack.com> @ #nlp-test

Expect Rapid Releases & Long-Term Support from John Snow Labs.