



# Machine Learning in Production Model Testing beyond Accuracy

= (Slicing, Capabilities, Invariants, Simulation, ...)

# Learning Goals

- Curate validation datasets for assessing model quality, covering subpopulations and capabilities as needed
- Explain the oracle problem and how it challenges testing of software and models
- Use invariants to check partial model properties with automated testing
- Select and deploy automated infrastructure to evaluate and monitor model quality

# Model Quality

## First Part: Measuring Prediction Accuracy

- the data scientist's perspective

## Second Part: What is Correctness Anyway?

- the role and lack of specifications, validation vs verification

## Third Part: Learning from Software Testing

- unit testing, test case curation, invariants, simulation (next lecture)

## Later: Testing in Production

- monitoring, A/B testing, canary releases (in 2 weeks)

# Curating Validation Data & Input Slicing



# Breakout Discussion

Write a few tests for the following program:

```
def nextDate(year: Int, month: Int, day: Int) = ...
```

A test may look like:

```
assert nextDate(2021, 2, 8) == (2021, 2, 9);
```

**As a group, discuss how you select tests. Discuss how many tests you need to feel confident.**

Post answer to #lecture tagging group members in Slack using template:

*Selection strategy: ...*

*Test quantity: ...*

# Defining Software Testing

- Program  $p$  with specification  $s$
- Test consists of
  - Controlled environment
  - Test call, test inputs
  - Expected behavior/output (oracle)

```
assertEquals(4, add(2, 2));  
assertEquals(??, factorPrime(15485863));
```

Testing is complete but unsound: Cannot guarantee the absence of bugs

# How to Create Test Cases?

```
def nextDate(year: Int, month: Int, day: Int) = ...
```



## Speaker notes

Can focus on specification (and concepts in the domain, such as leap days and month lengths) or can focus on implementation

Will not randomly sample from distribution of all days



# Software Test Case Design

**Opportunistic/exploratory testing:** Add some unit tests, without much planning

**Specification-based testing ("black box"):** Derive test cases from specifications

- Boundary value analysis
- Equivalence classes
- Combinatorial testing
- Random testing

**Structural testing ("white box"):** Derive test cases to cover implementation paths

- Line coverage, branch coverage
- Control-flow, data-flow testing, MCDC, ...

Test execution usually automated, but can be manual too; automated generation

☰ from specifications or code possible

# Example: Boundary Value Testing

Analyze the specification, not the implementation!

**Key Insight:** Errors often occur at the boundaries of a variable value

For each variable select (1) minimum, (2) min+1, (3) medium, (4) max-1, and (5) maximum; possibly also invalid values min-1, max+1

Example: `nextDate(2015, 6, 13) = (2015, 6, 14)`

- **Boundaries?**

# Example: Equivalence classes

Idea: Typically many values behave similarly, but some groups of values are different

Equivalence classes derived from specifications (e.g., cases, input ranges, error conditions, fault models)

Example `nextDate(2015, 6, 13)`

- leap years, month with 28/30/31 days, days 1-28, 29, 30, 31

Pick 1 value from each group, combine groups from all variables

# Exercise

```
/** Compute the price of a bus ride:  
 * - Children under 2 ride for free, children under 18 and  
 *   senior citizen over 65 pay half, all others pay the  
 *   full fare of $3.  
 * - On weekdays, between 7am and 9am and between 4pm and  
 *   7pm a peak surcharge of $1.5 is added.  
 * - Short trips under 5min during off-peak time are free.*/  
def busTicketPrice(age: Int,  
                    datetime: LocalDateTime,  
                    rideTime: Int)
```

*suggest test cases based on boundary value analysis and equivalence class testing*

# Selecting Validation Data for Model Quality?

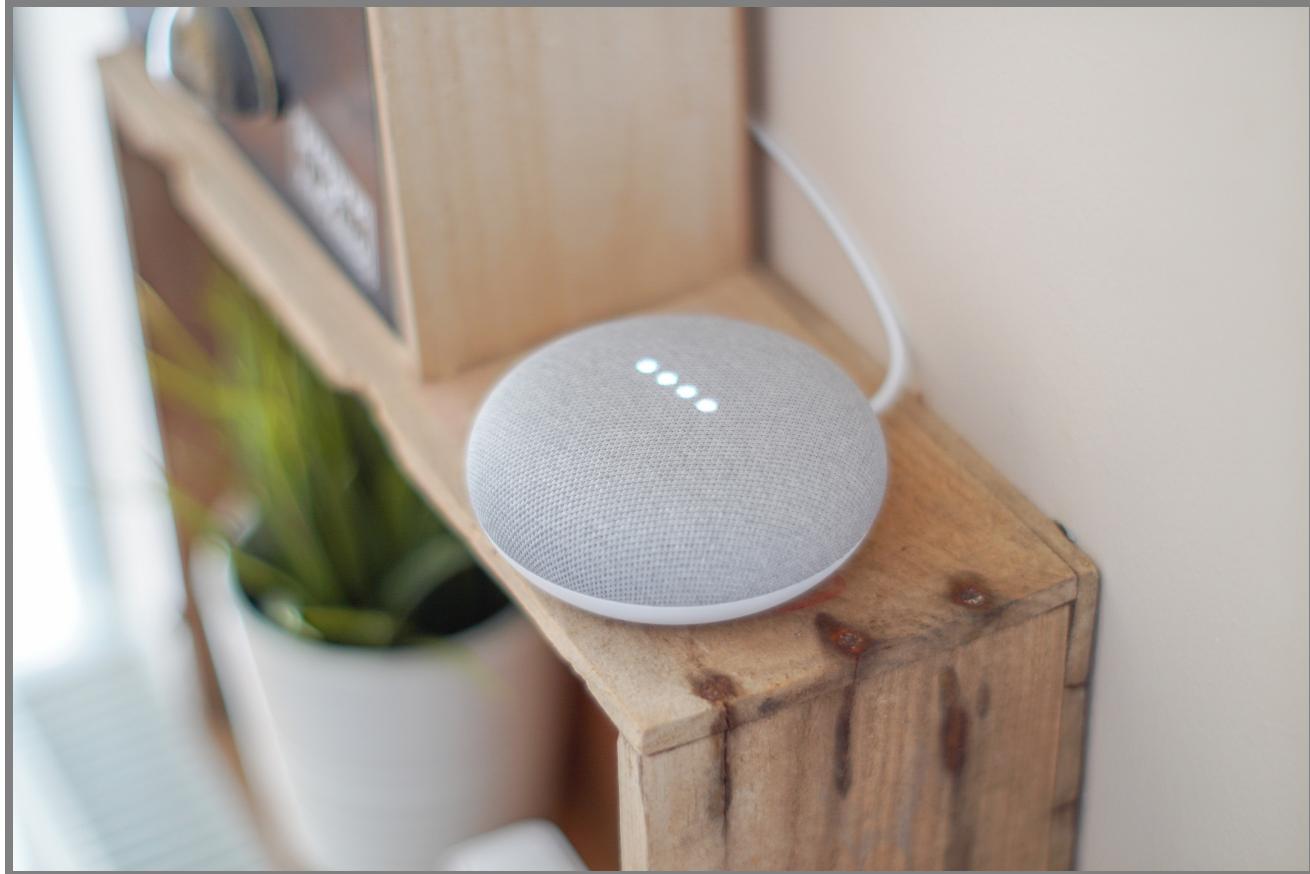


# Validation Data Representative?

- Validation data should reflect usage data
- Be aware of data drift (face recognition during pandemic, new patterns in credit card fraud detection)
- "*Out of distribution*" predictions often low quality (it may even be worth to detect out of distribution data in production, more later)

*(note, similar to requirements validation: did we hear all/representative stakeholders)*

# Not All Inputs are Equal



"Call mom" "What's the weather tomorrow?" "Add asafetida to my shopping list"

# Not All Inputs are Equal

*There Is a Racial Divide in Speech-Recognition Systems, Researchers Say: Technology from Amazon, Apple, Google, IBM and Microsoft misidentified 35 percent of words from people who were black. White people fared much better.* -- [NYTimes March 2020](#)



Chukwuemeka Afigbo  
@nke\_ise · [Follow](#)



If you have ever had a problem grasping the importance of diversity in tech and its impact on society, watch this video



9:48 AM · Aug 16, 2017



[Heart](#) 206K    [Reply](#)    [🔗 Copy link](#)

[Read 2.7K replies](#)

# Not All Inputs are Equal

*some random mistakes vs rare but biased mistakes?*

- A system to detect when somebody is at the door that never works for people under 5ft (1.52m)
- A spam filter that deletes alerts from banks

**Consider separate evaluations for important subpopulations;  
monitor mistakes in production**

# Identify Important Inputs

Curate Validation Data for Specific Problems and Subpopulations:

- *Regression testing*: Validation dataset for important inputs ("call mom") -- expect very high accuracy -- closest equivalent to **unit tests**
- *Uniformness/fairness testing*: Separate validation dataset for different subpopulations (e.g., accents) -- expect comparable accuracy
- *Setting goals*: Validation datasets for challenging cases or stretch goals -- accept lower accuracy

# Important Input Groups for Cancer Prognosis?



# Input Partitioning

- Guide testing by identifying groups and analyzing accuracy of subgroups
  - Often for fairness: gender, country, age groups, ...
  - Possibly based on business requirements or cost of mistakes
- Slice test data by population criteria, also evaluate interactions
- Identifies problems and plan mitigations, e.g., enhance with more data for subgroup or reduce confidence

Good reading: Barash, Guy, Eitan Farchi, Ilan Jayaraman, Orna Raz, Rachel Tzoref-Brill, and Marcel Zalmanovici. "Bridging the gap between ML solutions and their business requirements using feature

# Input Partitioning Example

DECADE	SUPPORT	ACC
1910s	38	78.94
1930s	338	87.87
1990s	3007	90.95
2000s	6192	91.40

Input divided by movie age. Notice low accuracy, but also low support (i.e., little validation data), for old movies.

MAIN_GENRE	RAT_CAT	LEN_CAT	SUPPORT	ACC
Mystery	OK	long	11	72.72
Fantasy	OK	short	36	77.77
Crime	OK	long	100	81.00
Comedy	GOOD	long	55	96.36

Input divided by genre, rating, and length. Accuracy differs, but also amount of test data used ("support") differs, highlighting low confidence areas.

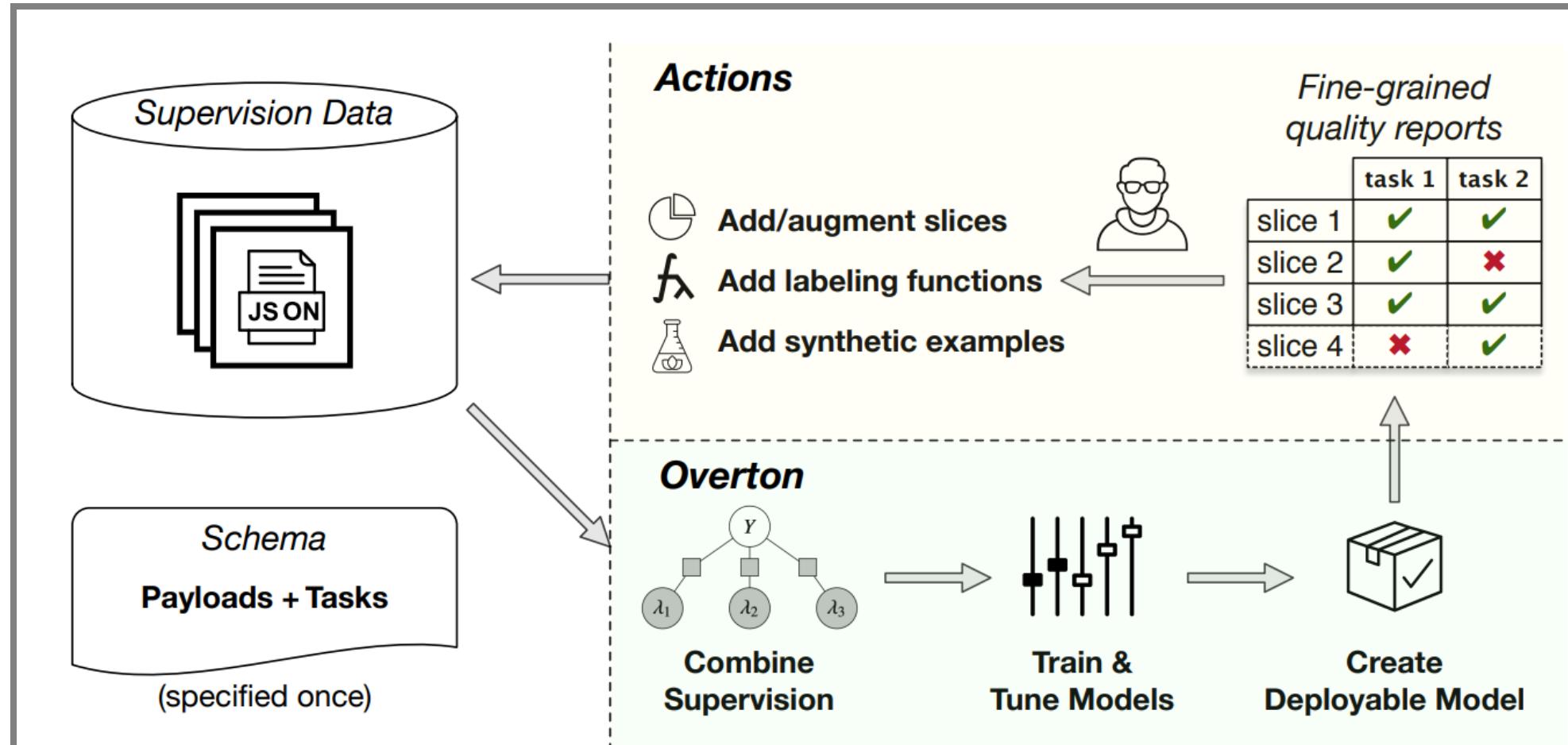
Source: Barash, Guy, et al. "Bridging the gap between ML solutions and their business requirements using feature interactions." In Proc. FSE, 2019.

# Input Partitioning Discussion

How to slice evaluation data for cancer prognosis?



# Example: Model Impr. at Apple (Overton)



Ré, Christopher, Feng Niu, Pallavi Gudipati, and Charles Srisuwananukorn. "[Overton: A Data System for Monitoring and Improving Machine-Learned Products](#)." arXiv preprint arXiv:1909.05372 (2019).

# Example: Model Improvement at Apple (Overton)

- Focus engineers on creating training and validation data, not on model search (AutoML)
- Flexible infrastructure to slice telemetry data to identify underperforming subpopulations -> focus on creating better training data (better, more labels, in semi-supervised learning setting)

# Testing Model Capabilities



Further reading: Christian Kaestner. [Rediscovering Unit Testing: Testing Capabilities of ML Models](#).  
Toward Data Science, 2021.

# Testing Capabilities

Are there "concepts" or "capabilities" the model should learn?

Example capabilities of sentiment analysis:

- Handle *negation*
- Robustness to *typos*
- Ignore synonyms and abbreviations
- Person and location names are irrelevant
- Ignore gender
- ...

For each capability create specific test set (multiple examples)

Ribeiro, Marco Tulio, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. "[Beyond Accuracy: Behavioral Testing of NLP Models with CheckList](#)." In Proceedings ACL, p. 4902–4912. (2020).

# Testing Capabilities

Robust.	<b>INV:</b> Add randomly generated URLs and handles to tweets	9.6	13.4	24.8	11.4	7.4	@JetBlue that selfie was extreme. <a href="#">@pi9QDK</a> INV @united stuck because staff took a break? Not happy 1K.... <a href="https://t.co/PWK1jb">https://t.co/PWK1jb</a> INV
	<b>INV:</b> Swap one character with its neighbor (typo)	5.6	10.2	10.4	5.2	3.8	<a href="#">@JetBlue</a> → <a href="#">@JeBtblue</a> I cri INV @SouthwestAir no <a href="#">thanks</a> → <a href="#">thakns</a> INV
NER	<b>INV:</b> Switching locations should not change predictions	7.0	20.8	14.8	7.6	6.4	@JetBlue I want you guys to be the first to fly to # <a href="#">Cuba</a> → <a href="#">Canada</a> ... INV @VirginAmerica I miss the #nerdbird in <a href="#">San Jose</a> → <a href="#">Denver</a> INV
	<b>INV:</b> Switching person names should not change predictions	2.4	15.1	9.1	6.6	2.4	...Airport agents were horrendous. <a href="#">Sharon</a> → <a href="#">Erin</a> was your saviour INV @united 8602947, <a href="#">Jon</a> → <a href="#">Sean</a> at <a href="http://t.co/58tuTgli0D">http://t.co/58tuTgli0D</a> , thanks. INV

From: Ribeiro, Marco Tulio, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. "Beyond Accuracy: Behavioral Testing of NLP Models with Checklist." In Proceedings ACL, p. 4902–4912. (2020).

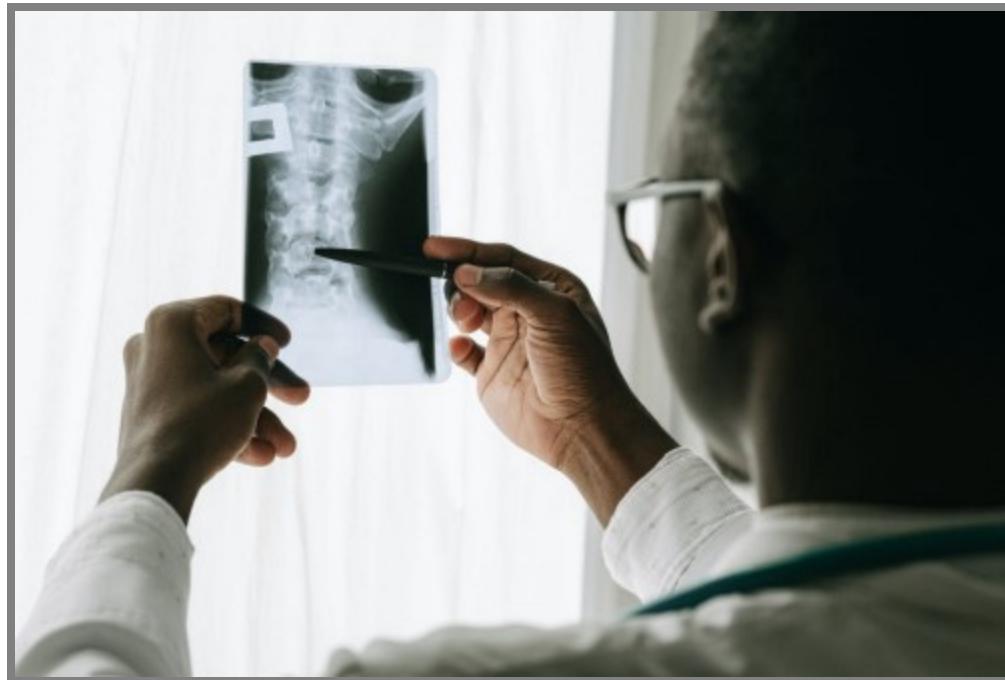
# Testing Capabilities

Negation	<b>MFT:</b> Negated negative should be positive or neutral	18.8	54.2	29.4	13.2	2.6	The food is not poor. pos or neutral It isn't a lousy customer service. pos or neutral
	<b>MFT:</b> Negated neutral should still be neutral	40.4	39.6	74.2	98.4	95.4	This aircraft is not private. neutral This is not an international flight. neutral
	<b>MFT:</b> Negation of negative at the end, should be pos. or neut.	100.0	90.4	100.0	84.8	7.2	I thought the plane would be awful, but it wasn't. pos or neutral I thought I would dislike that plane, but I didn't. pos or neutral
	<b>MFT:</b> Negated positive with neutral content in the middle	98.4	100.0	100.0	74.0	30.2	I wouldn't say, given it's a Tuesday, that this pilot was great. neg I don't think, given my history with airplanes, that this is an amazing staff. neg
	<b>MFT:</b> Author sentiment is more important than of others	45.4	62.4	68.0	38.8	30.0	Some people think you are excellent, but I think you are nasty. neg Some people hate you, but I think you are exceptional. pos

From: Ribeiro, Marco Tulio, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. "[Beyond Accuracy: Behavioral Testing of NLP Models with CheckList](#)." In Proceedings ACL, p. 4902–4912. (2020).

# Examples of Capabilities

What could be capabilities of the cancer classifier?



# Capabilities vs Specifications vs Slicing

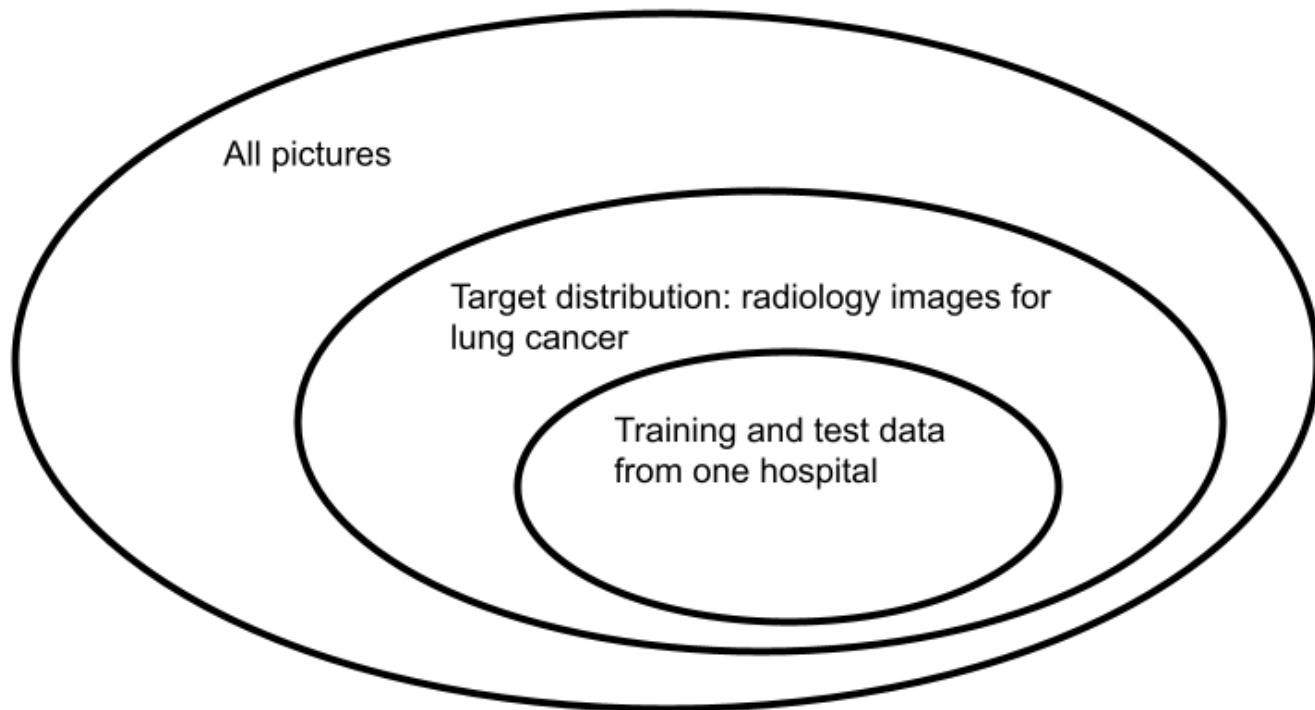


# Capabilities vs Specifications vs Slicing

Capabilities are partial specifications of expected behavior (not expected to always hold)

Some capabilities correspond to slices of existing test data, for others we may need to create new data

# Recall: Is it fair to expect generalization beyond training distribution?



*Shall a cancer detector generalize to other hospitals? Shall image captioning generalize to describing pictures of star formations?*

## Speaker notes

We wouldn't test a first year elementary school student on high-school math. This would be "out of the training distribution"



# Recall: Shortcut Learning

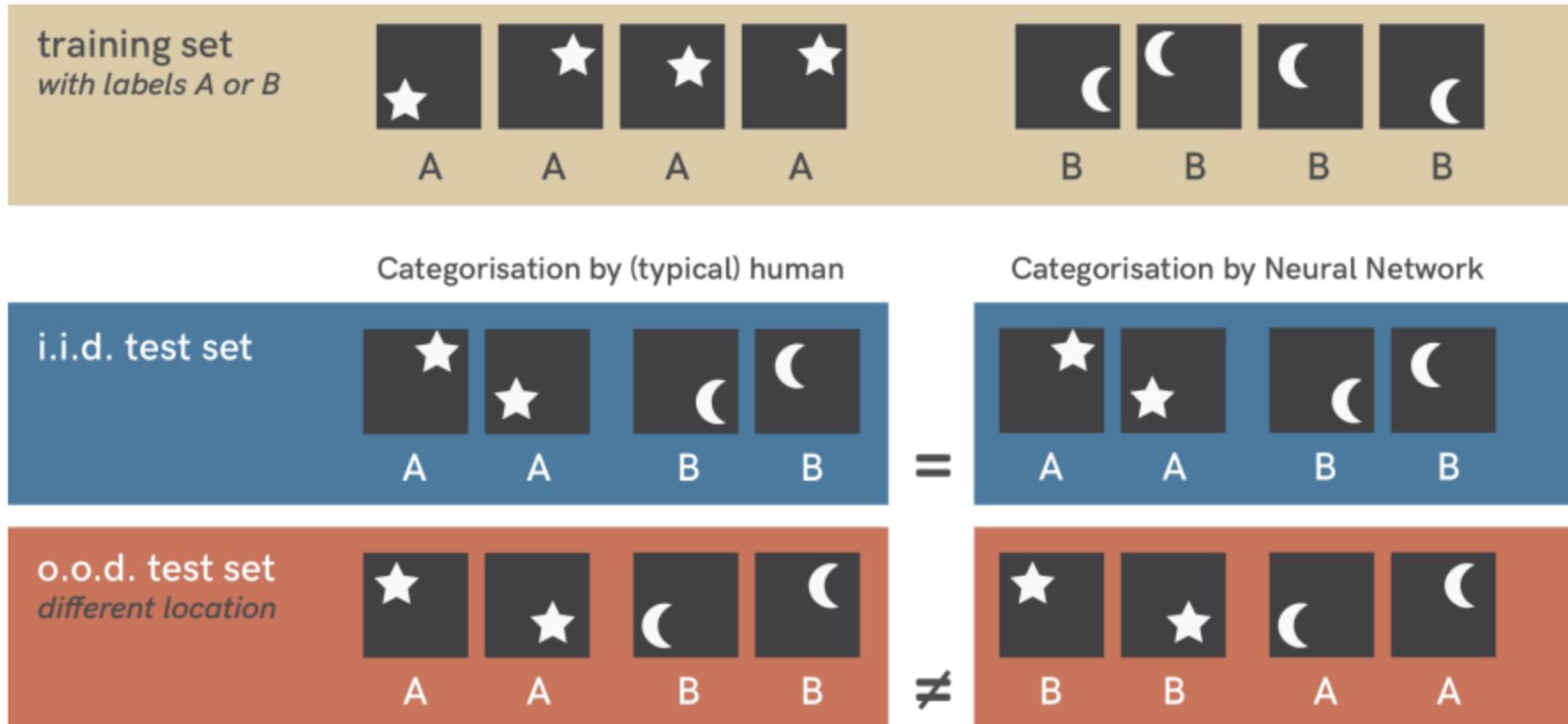


Figure from: Geirhos, Robert, et al. "[Shortcut learning in deep neural networks](#)." Nature Machine Intelligence 2, no. 11 (2020): 665-673.

# More Shortcut Learning :)



(A) **Cow: 0.99**, Pasture: 0.99, Grass: 0.99, No Person: 0.98, Mammal: 0.98



(B) No Person: 0.99, Water: 0.98, Beach: 0.97, Outdoors: 0.97, Seashore: 0.97



(C) No Person: 0.97, **Mammal: 0.96**, Water: 0.94, Beach: 0.94, Two: 0.94

Figure from Beery, Sara, Grant Van Horn, and Pietro Perona. "Recognition in terra incognita." In Proceedings of the European Conference on Computer Vision (ECCV), pp. 456–473. 2018.

# Generalization beyond Training Distribution?

- Typically training and validation data from same distribution (i.i.d. assumption!)
- Many models can achieve similar accuracy
- Models that learn "right" abstractions possibly indistinguishable from models that use shortcuts
  - see tank detection example
  - Can we guide the model towards "right" abstractions?
- Some models generalize better to other distributions not used in training
  - e.g., cancer images from other hospitals, from other populations
  - Drift and attacks, ...

# Hypothesis: Testing Capabilities may help with Generalization

- Capabilities are "partial specifications", given beyond training data
- Encode domain knowledge of the problem
  - Capabilities are inherently domain specific
  - Curate capability-specific test data for a problem
- Testing for capabilities helps to distinguish models that use intended abstractions
- May help find models that generalize better

See discussion in D'Amour, Alexander, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen et al. "[Underspecification presents challenges for credibility in modern machine learning.](#)" arXiv preprint arXiv:2011.03395 (2020).

# Strategies for identifying capabilities

- Analyze common mistakes (e.g., classify past mistakes in cancer prognosis)
- Use existing knowledge about the problem (e.g., linguistics theories)
- Observe humans (e.g., how do radiologists look for cancer)
- Derive from requirements (e.g., fairness)
- Causal discovery from observational data?

# Examples of Capabilities

What could be capabilities of image captioning system?



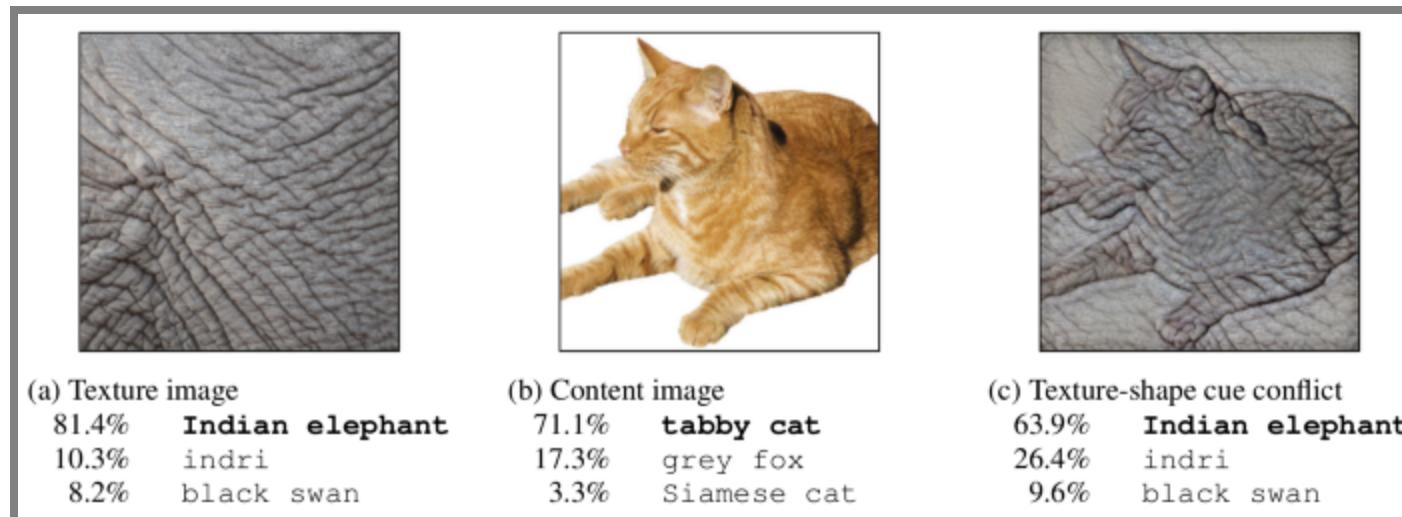
# Generating Test Data for Capabilities

## Idea 1: Domain-specific generators

Testing *negation* in sentiment analysis with template:

I {NEGATION} {POS\_VERB} the {THING}.

Testing texture vs shape priority with artificial generated images:



# Generating Test Data for Capabilities

## Idea 2: Mutating existing inputs

Testing *synonyms* in sentiment analysis by replacing words with synonyms, keeping label

Testing *robust against noise and distraction* add and false is not true or random URLs to text

# Generating Test Data for Capabilities

## Idea 3: Crowd-sourcing test creation

Testing *sarcasm* in sentiment analysis: Ask humans to minimally change text to flip sentiment with sarcasm

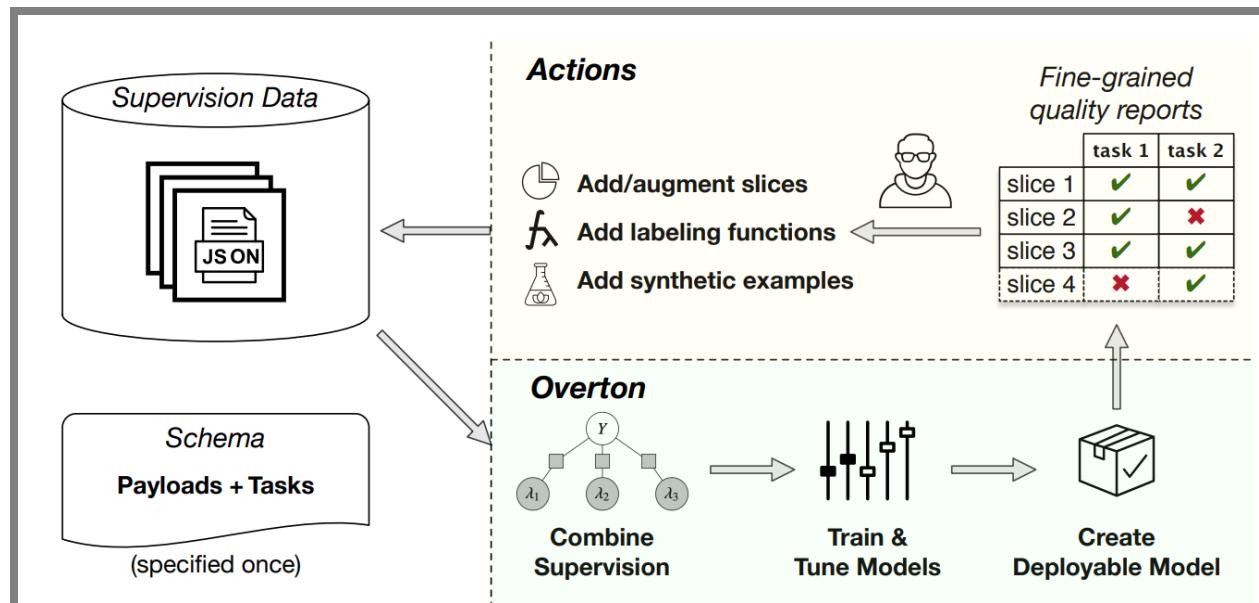
Testing *background* in object detection: Ask humans to take pictures of specific objects with unusual backgrounds

Recasting fact as hoped for	The world of Atlantis, hidden beneath the earth's core, is fantastic The world of Atlantis, hidden beneath the earth's core is <b>supposed</b> to be fantastic
Suggesting sarcasm	thoroughly captivating <b>thriller-drama, taking a deep and realistic</b> view thoroughly mind numbing " <b>thriller-drama”, taking a “deep” and “realistic” (who are they kidding?)</b> view
Inserting modifiers	The presentation of simply Atlantis' landscape and setting The presentation of Atlantis' <b>predictable</b> landscape and setting

# Generating Test Data for Capabilities

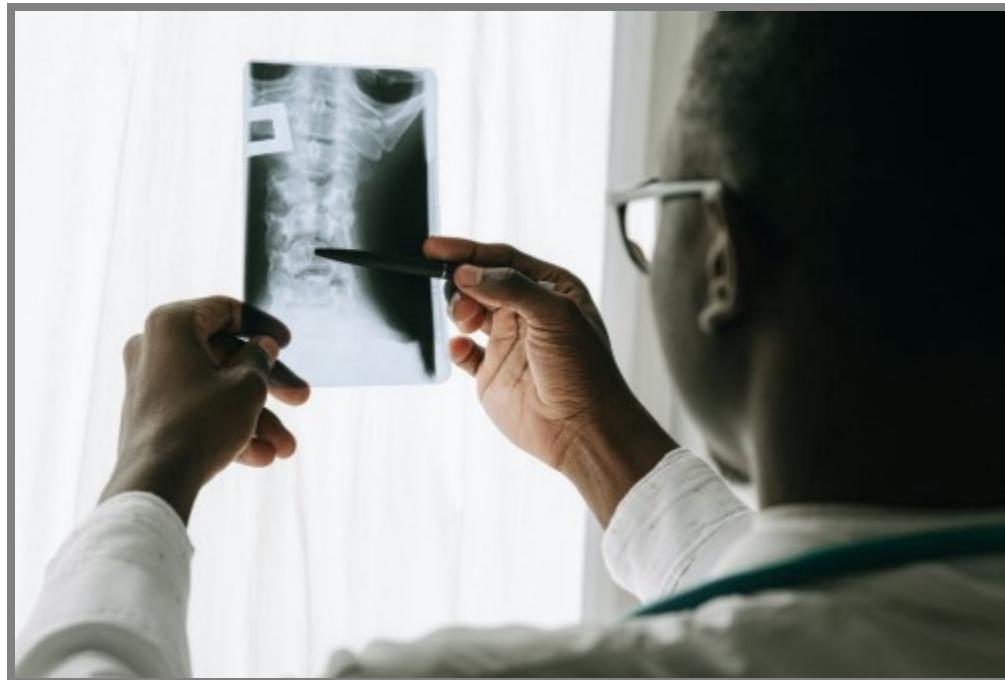
## Idea 4: Slicing test data

Testing *negation* in sentiment analysis by finding sentences containing 'not'



# Examples of Capabilities

How to generate test data for capabilities of the cancer classifier?



# Testing vs Training Capabilities

- Dual insight for testing and training
- Strategies for curating test data can also help select training data
- Generate capability-specific training data to guide training (data augmentation)

Further reading on using domain knowledge during training: Von Rueden, Laura, Sebastian Mayer, Jochen Garcke, Christian Bauckhage, and Jannis Schuecker. "Informed machine learning—towards a taxonomy of explicit integration of knowledge into machine learning." Learning 18 (2019): 19-20.

# Preliminary Summary: Specification-Based Testing Techniques as Inspiration

- Boundary value analysis
- Partition testing & equivalence classes
- Combinatorial testing
- Decision tables

Use to identify datasets for **subpopulations** and **capabilities**, not individual tests.

# On Terminology



- Test data curation is emerging as a very recent concept for testing ML components
- No consistent terminology
  - "Testing capabilities" in checklist paper
  - "Stress testing" in some others (but stress testing has a very different meaning in software testing: robustness to overload)
- Software engineering concepts translate, but names not adopted in ML community
  - specification-based testing, black-box testing
  - equivalence class testing, boundary-value analysis

# Automated (Random) Testing and Invariants

(if it wasn't for that darn oracle problem)

# Random Test Input Generation is Easy

```
@Test  
void testNextDate() {  
    nextDate(488867101, 1448338253, -997372169)  
    nextDate(2105943235, 1952752454, 302127018)  
    nextDate(1710531330, -127789508, 1325394033)  
    nextDate(-1512900479, -439066240, 889256112)  
    nextDate(1853057333, 1794684858, 1709074700)  
    nextDate(-1421091610, 151976321, 1490975862)  
    nextDate(-2002947810, 680830113, -1482415172)  
    nextDate(-1907427993, 1003016151, -2120265967)  
}
```

But is it useful?

# Cancer in Random Image?

# Randomly Generating "Realistic" Inputs is Possible

```
@Test  
void testNextDate() {  
    nextDate(2010, 8, 20)  
    nextDate(2024, 7, 15)  
    nextDate(2011, 10, 27)  
    nextDate(2024, 5, 4)  
    nextDate(2013, 8, 27)  
    nextDate(2010, 2, 30)  
}
```

But how do we know whether the computation is correct?

# Automated Model Validation Data Generation?

```
@Test  
void testCancerPrediction() {  
    cancerModel.predict(generateRandomImage())  
    cancerModel.predict(generateRandomImage())  
    cancerModel.predict(generateRandomImage())  
}
```

- **Realistic inputs?**
- **But how do we get labels?**

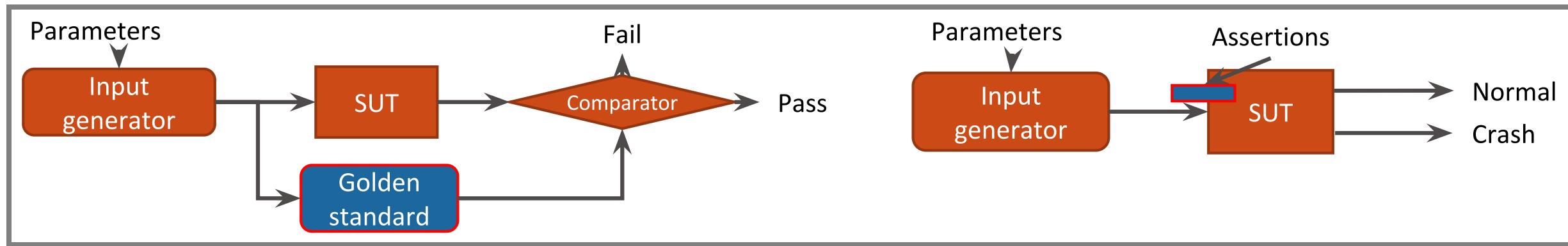
# The Oracle Problem

*How do we know the expected output of a test?*

```
assertEquals(??, factorPrime(15485863));
```

# Test Case Generation & The Oracle Problem

- Manually construct input-output pairs (does not scale, cannot automate)
- Comparison against gold standard (e.g., alternative implementation, executable specification)
- Checking of global properties only -- crashes, buffer overflows, code injections
- Manually written assertions -- partial specifications checked at runtime



# Manually constructing outputs

```
@Test  
void testNextDate() {  
    assert nextDate(2010, 8, 20) == (2010, 8, 21);  
    assert nextDate(2024, 7, 15) == (2024, 7, 16);  
    assert nextDate(2010, 2, 30) throws InvalidInputException;  
}
```

```
@Test  
void testCancerPrediction() {  
    assert cancerModel.predict(loadImage("random1.jpg")) == true  
    assert cancerModel.predict(loadImage("random2.jpg")) == true  
    assert cancerModel.predict(loadImage("random3.jpg")) == false  
}
```

# Compare against reference implementation

## assuming we have a correct implementation

```
@Test  
void testNextDate() {  
    assert nextDate(2010, 8, 20) == referenceLib.nextDate(2010,  
    assert nextDate(2024, 7, 15) == referenceLib.nextDate(2024,  
    assert nextDate(2010, 2, 30) == referenceLib.nextDate(2010,  
}
```

```
@Test  
void testCancerPrediction() {  
    assert cancerModel.predict(loadImage("random1.jpg")) == ???;  
}
```

# Checking global specifications

Ensure, no computation crashes

```
@Test  
void testNextDate() {  
    nextDate(2010, 8, 20)  
    nextDate(2024, 7, 15)  
    nextDate(2010, 2, 30)  
}
```

```
@Test  
void testCancerPrediction() {  
    cancerModel.predict(generateRandomImage())  
    cancerModel.predict(generateRandomImage())
```

# Invariants as partial specification

```
class Stack {  
    int size = 0;  
    int MAX_SIZE = 100;  
    String[] data = new String[MAX_SIZE];  
    // class invariant checked before and after every method  
    private void check() {  
        assert(size>=0 && size<=MAX_SIZE);  
    }  
    public void push(String v) {  
        check();  
        if (size<MAX_SIZE)  
    }
```

# Automated Testing / Test Case Generation / Fuzzing

- Many techniques to generate test cases
  - Dumb fuzzing: generate random inputs
  - Smart fuzzing (e.g., symbolic execution, coverage guided fuzzing): generate inputs to maximally cover the implementation
  - Program analysis to understand the shape of inputs, learning from existing tests
  - Minimizing redundant tests
  - Abstracting/simulating/mock the environment
- ≡ • Typically looking for crashing bugs or assertion violations

# Test Generation (Symbolic Execution)

Code:

```
void foo(a, b, c) {  
    int x=0, y=0, z=0;  
    if (a) x=-2;  
    if (b<5) {  
        if (!a && c) y=1;  
        z=2;  
    }  
    assert(x+y+z!=3)  
}
```

Paths:

- :  $x = -2, y = 0, z = 2$
- $a \wedge x = -2, y = 0, z = 0$
- $a \wedge 0 \leq b < 5 \wedge z = 2$
- $\neg a \wedge 0 \leq z = 0, z = 2$
- $\neg a \wedge 0 \leq b < 5 \wedge \neg(a \wedge c)$
- $\neg a \wedge 0 \leq b < 5 \wedge \neg(a \wedge c) \wedge \neg a \wedge \neg(b < 5)$

## Speaker notes

example source: <http://web.cs.iastate.edu/~weile/cs641/9.SymbolicExecution.pdf>



# Generating Inputs for ML Problems

- Completely random data generation (uniform sampling from each feature's domain)
- Using knowledge about feature distributions (sample from each feature's distribution)
- Knowledge about dependencies among features and whole population distribution (e.g., model with probabilistic programming language)
- Mutate from existing inputs (e.g., small random modifications to select features)
- Generate "fake data" with Generative Adversarial Networks

# ML Models = Untestable Software?

```
@Test  
void testCancerPrediction() {  
    cancerModel.predict(generateRandomImage())  
}
```

- Manually construct input-output pairs (does not scale, cannot automate)
  - **too expensive at scale**
- Comparison against gold standard (e.g., alternative implementation, executable specification)
  - **no specification, usually no other "correct" model**
  - comparing different techniques useful? (see ensemble learning)
  - semi-supervised learning as approximation?
- Checking of global properties only -- crashes, buffer overflows, code injections - ??
- Manually written assertions -- partial specifications checked at runtime - ??

# Invariants in Machine Learned Models (Metamorphic Testing)

Exploit relationships between inputs

- If two inputs differ only in X -> output should be the same
- If inputs differ in Y output should be flipped
- If inputs differ only in feature F, prediction for input with higher F should be higher
- ...

# Invariants in Machine Learned Models?



# Some Capabilities are Invariants

Some capability tests can be expressed as invariants and automatically encoded as transformations to existing test data

- Negation should flip sentiment analysis result
- Typos should not affect sentiment analysis result
- Changes to locations or names should not affect sentiment analysis results

Robust.	<i>INV:</i> Add randomly generated URLs and handles to tweets	9.6	13.4	24.8	11.4	7.4	@JetBlue that selfie was extreme. @pi9QDK INV @united stuck because staff took a break? Not happy 1K.... <a href="https://t.co/PWK1jb">https://t.co/PWK1jb</a> INV
	<i>INV:</i> Swap one character with its neighbor (typo)	5.6	10.2	10.4	5.2	3.8	@JetBlue → @JeBtue I cri INV @SouthwestAir no thanks → thakns INV
NER	<i>INV:</i> Switching locations should not change predictions	7.0	20.8	14.8	7.6	6.4	@JetBlue I want you guys to be the first to fly to # Cuba → Canada... INV @VirginAmerica I miss the #nerdbird in San Jose → Denver INV
	<i>INV:</i> Switching person names should not change predictions	2.4	15.1	9.1	6.6	2.4	...Airport agents were horrendous. Sharon → Erin was your saviour INV @united 8602947, Jon → Sean at <a href="http://t.co/58tuTgli0D">http://t.co/58tuTgli0D</a> , thanks. INV

# Examples of Invariants

- Credit rating should not depend on gender:
  - $\forall x. f(x[\text{gender} \leftarrow \text{male}]) = f(x[\text{gender} \leftarrow \text{female}])$
- Synonyms should not change the sentiment:
  - $\forall x. f(x[\text{synonym} \leftarrow \text{synonym}']) = f(x)$
- Negation should swap meaning:
  - $\forall x. \neg f(x) \equiv f(\text{not } x)$
- Robustness around training data:
  - $\forall x. f(\text{replace}(x, " is ", " is not "))$
- Low credit scores should never get a loan (sufficient conditions for classification, "anchors"):
  - $\forall x. x.\text{score} < 649 \Rightarrow \neg f(x)$

Identifying invariants requires domain knowledge of the problem!

# Metamorphic Testing

Formal description of relationships among inputs and outputs  
(*Metamorphic Relations*)

In general, for a model and inputs define two functions to transform inputs and outputs and such that:<sup>x</sup>

$$g_I \quad g_O$$

$$\forall x. f(g_I(x)) = g_O(f(x))$$

e.g. and

$$g_I(x) \neq \text{replace}(x, " \text{is} ", " \text{is not} ")$$

# On Testing with Invariants/Assertions

- Defining good metamorphic relations requires knowledge of the problem domain
- Good metamorphic relations focus on parts of the system
- Invariants usually cover only one aspect of correctness -- maybe capabilities
- Invariants and near-invariants can be mined automatically from sample data (see *specification mining* and *anchors*)

Further reading:

- Segura, Sergio, Gordon Fraser, Ana B. Sanchez, and Antonio Ruiz-Cortés. "[A survey on metamorphic testing.](#)" IEEE Transactions on software engineering 42, no. 9 (2016): 805-824.
- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "[Anchors: High-precision model-agnostic explanations.](#)" In Thirty-Second AAAI Conference on Artificial Intelligence. 2018.

# Invariant Checking aligns with Requirements Validation

# Approaches for Checking in Variants

- Generating test data (random, distributions) usually easy
- Transformations of existing test data
- Adversarial learning: For many techniques gradient-based techniques to search for invariant violations -- that's roughly analogous to symbolic execution in SE
- Early work on formally verifying invariants for certain models (e.g., small deep neural networks)

Further readings: Singh, Gagandeep, Timon Gehr, Markus Püschel, and Martin Vechev. "[An abstract domain for certifying neural networks.](#)" Proceedings of the ACM on Programming Languages 3, no. POPL (2019): 1-30.

# Using Invariant Violations

# Simulation-Based Testing



# One More Thing: Simulation-Based Testing

In some cases it is easy to go from outputs to inputs:

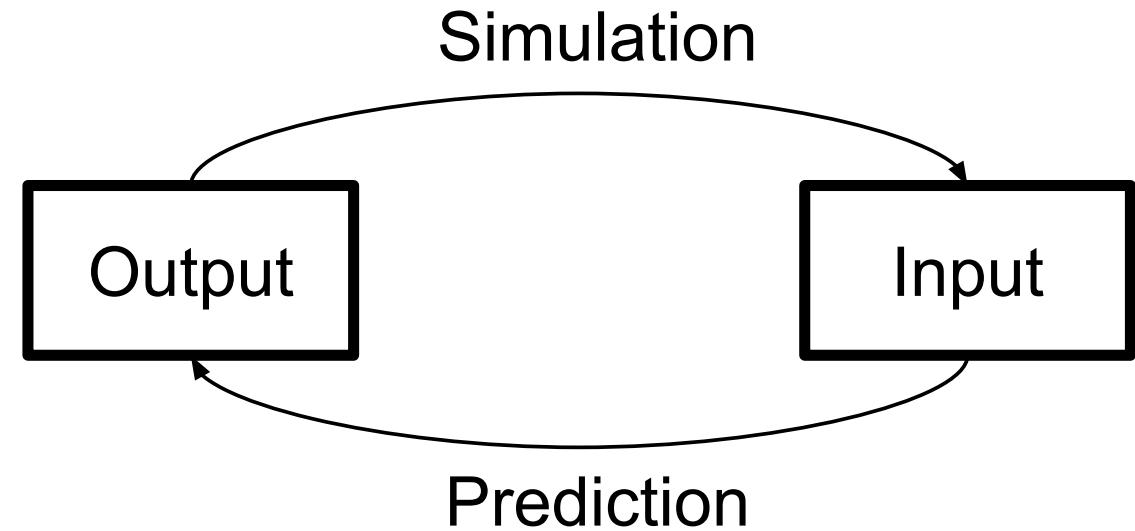
```
assertEquals(??, factorPrime(15485862));
```

```
randomNumbers = [2, 3, 7, 7, 52673]
assertEquals(randomNumbers,
    factorPrime(multiply(randomNumbers)));
```

Similar idea in machine-learning problems?

# One More Thing: Simulation-Based Testing

- Derive input-output pairs from simulation, esp. in vision systems
- Example: Vision for self-driving cars:
  - Render scene -> add noise -> recognize -> compare recognized result with simulator state
- Quality depends on quality of simulator:
  - examples: render picture/video, synthesize speech, ...
  - Less suitable where input-output relationship unknown, e.g., cancer prognosis, housing price prediction



Further readings: Zhang, Mengshi, Yuqun Zhang, Lingming Zhang, Cong Liu, and Sarfraz Khurshid. "DeepRoad: GAN-based metamorphic testing and input validation framework for autonomous driving systems." In Proc. ASE. 2018.

# Preliminary Summary: Invariants and Generation

- Generating sample inputs is easy, but knowing corresponding outputs is not (oracle problem)
- Crashing bugs are not a concern
- Invariants + generated data can check capabilities or properties (metamorphic testing)
  - Inputs can be generated realistically or to find violations (adversarial learning)
- If inputs can be computed from outputs, tests can be automated (simulation-based testing)

# On Terminology



**Metamorphic testing** is an academic software engineering term that's not common in ML literature, it generalizes many concepts regularly reinvented

Much of the security, safety and robustness literature in ML focuses on invariants

# Other Testing Concepts

# Test Coverage

Packages

All

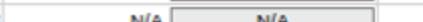
- [net.sourceforge.cobertura.ant](#)
- [net.sourceforge.cobertura.check](#)
- [net.sourceforge.cobertura.coveragedata](#)
- [net.sourceforge.cobertura.instrument](#)
- [net.sourceforge.cobertura.merge](#)
- [net.sourceforge.cobertura.reporting](#)
- [net.sourceforge.cobertura.reporting.html](#)
- [net.sourceforge.cobertura.reporting.html](#)
- [net.sourceforge.cobertura.reporting.xml](#)
- [net.sourceforge.cobertura.util](#)

All Packages

Classes

- [AntUtil \(88%\)](#)
- [Archive \(100%\)](#)
- [ArchiveUtil \(80%\)](#)
- [BranchCoverageData \(N/A\)](#)
- [CheckTask \(0%\)](#)
- [ClassData \(N/A\)](#)
- [ClassInstrumenter \(94%\)](#)
- [ClassPattern \(100%\)](#)
- [CoberturaFile \(73%\)](#)
- [CommandLineBuilder \(96%\)](#)
- [CommonMatchingTask \(88%\)](#)
- [ComplexityCalculator \(100%\)](#)
- [ConfigurationUtil \(50%\)](#)
- [CopyFiles \(87%\)](#)
- [CoverageData \(N/A\)](#)
- [CoverageDataContainer \(N/A\)](#)
- [CoverageDataFileHandler \(N/A\)](#)

Coverage Report - All Packages

Package	# Classes	Line Coverage	Branch Coverage	Complexity
All Packages	55	75% 	64% 	2.319
<a href="#">net.sourceforge.cobertura.ant</a>	11	52% 	43% 	1.848
<a href="#">net.sourceforge.cobertura.check</a>	3	0% 	0% 	2.429
<a href="#">net.sourceforge.cobertura.coveragedata</a>	13	N/A 	N/A 	2.277
<a href="#">net.sourceforge.cobertura.instrument</a>	10	90% 	75% 	1.854
<a href="#">net.sourceforge.cobertura.merge</a>	1	86% 	88% 	5.5
<a href="#">net.sourceforge.cobertura.reporting</a>	3	87% 	80% 	2.882
<a href="#">net.sourceforge.cobertura.reporting.html</a>	4	91% 	77% 	4.444
<a href="#">net.sourceforge.cobertura.reporting.html.files</a>	1	87% 	62% 	4.5
<a href="#">net.sourceforge.cobertura.reporting.xml</a>	1	100% 	95% 	1.524
<a href="#">net.sourceforge.cobertura.util</a>	9	60% 	69% 	2.892
<a href="#">someotherpackage</a>	1	83% 	N/A 	1.2

Report generated by Cobertura 1.9 on 6/9/07 12:37 AM.

# Example: Structural testing

```
int divide(int A, int B) {  
    if (A==0)  
        return 0;  
    if (B==0)  
        return -1;  
    return A / B;  
}
```

*minimum set of test cases to cover all lines? all decisions? all path?*

Packages	Coverage Report - All Packages				
	Package	# Classes	Line Coverage	Branch Coverage	Complexity
All	All Packages	55	75% <div style="width: 75%; background-color: green; height: 10px;"></div> 1625/2179	64% <div style="width: 64%; background-color: green; height: 10px;"></div> 472/738	2.319
<a href="#">net.sourceforge.cobertura.ant</a>	net.sourceforge.cobertura.ant	11	52% <div style="width: 52%; background-color: green; height: 10px;"></div> 170/330	43% <div style="width: 43%; background-color: green; height: 10px;"></div> 40/94	1.848
<a href="#">net.sourceforge.cobertura.check</a>	net.sourceforge.cobertura.check	3	0% <div style="width: 0%; background-color: red; height: 10px;"></div> 0/150	0% <div style="width: 0%; background-color: red; height: 10px;"></div> 0/76	2.429
<a href="#">net.sourceforge.cobertura.coveragedata</a>	net.sourceforge.cobertura.coveragedata	13	N/A <div style="width: 0%; background-color: grey; height: 10px;"></div> N/A	N/A <div style="width: 0%; background-color: grey; height: 10px;"></div> N/A	2.277
<a href="#">net.sourceforge.cobertura.instrument</a>	net.sourceforge.cobertura.instrument	10	90% <div style="width: 90%; background-color: green; height: 10px;"></div> 460/510	75% <div style="width: 75%; background-color: green; height: 10px;"></div> 123/164	1.854
<a href="#">net.sourceforge.cobertura.merge</a>	net.sourceforge.cobertura.merge	1	86% <div style="width: 86%; background-color: green; height: 10px;"></div> 30/35	88% <div style="width: 88%; background-color: green; height: 10px;"></div> 14/16	5.5

# Defining Structural Testing ("white box")

- Test case creation is driven by the implementation, not the specification
- Typically aiming to increase coverage of lines, decisions, etc
- Automated test generation often driven by maximizing coverage (for finding crashing bugs)

# Whitebox Analysis in ML

- Several coverage metrics have been proposed
  - All path of a decision tree?
  - All neurons activated at least once in a DNN? (several papers "neuron coverage")
  - Linear regression models??
- Often create artificial inputs, not realistic for distribution
- Unclear whether those are useful
- Adversarial learning techniques usually more efficient at finding invariant violations

# Regression Testing

- Whenever bug detected and fixed, add a test case
- Make sure the bug is not reintroduced later
- Execute test suite after changes to detect regressions
  - Ideally automatically with continuous integration tools
- Maps well to curating test sets for important populations in ML

# Mutation Analysis

- Start with program and passing test suite
- Automatically insert small modifications ("mutants") in the source code
  - $a+b \rightarrow a-b$
  - $a < b \rightarrow a \leq b$
  - ...
- Can program detect modifications ("kill the mutant")?
- Better test suites detect more modifications ("mutation score")

```
int divide(int A, int B) {  
    if (A==0)          // A!=0, A<0, B==0  
        return 0;      // 1, -1
```

# Mutation Analysis

- Some papers exist, but strategy unclear
- Mutating model parameters? Mutating hyperparameters? Mutating inputs?
- What's considered as killing a mutant, if we don't have specifications?
- Still unclear application...

# Continuous Integration for Model Quality

# Continuous Integration

The screenshot shows a web browser window displaying the Travis CI interface for a repository named "wyvernlang / wyvern". The build number is #17, and the status is "passing". The build was authored and committed by "potanin". The build duration was 16 seconds, and it finished 3 days ago. The build ran on a Linux worker. A message at the bottom of the log indicates that the job ran on legacy infrastructure and suggests upgrading.

Build #17 - wyvernlang

https://travis-ci.org/wyvernlang/wyvern/builds/79099642

Travis CI Blog Status Help Jonathan Aldrich

Search all repositories

wyvernlang / wyvern build passing

My Repositories +

Current Branches Build History Pull Requests > Build #17 Settings

wyvernlang/wyvern # 17

Duration: 16 sec Finished: 3 days ago

SimpleWyvern-devel Asserting false (works on Linux, so its OK). # 17 passed Commit fd7be1c Compare 0e2af1f..fd7be1c ran for 16 sec 3 days ago

potanin authored and committed

This job ran on our legacy infrastructure. Please read [our docs on how to upgrade](#)

Remove Log Download Log

Using worker: worker-linux-027f0490-1.bb.travis-ci.org:travis-linux-2

Build system information

system\_info

# Continuous Integration for Model Quality?



# Continuous Integration for Model Quality

- Testing script
  - Existing model: Automatically evaluate model on labeled training set; multiple separate evaluation sets possible, e.g., for slicing, regressions
  - Training model: Automatically train and evaluate model, possibly using cross-validation; many ML libraries provide built-in support
  - Report accuracy, recall, etc. in console output or log files
  - May deploy learning and evaluation tasks to cloud services
  - Optionally: Fail test below bound (e.g., accuracy <.9; accuracy < last accuracy)
- Version control test data, model and test scripts, ideally also learning data and learning code (feature extraction, modeling, ...)
- Continuous integration tool can trigger test script and parse output, plot for comparisons (e.g., similar to performance tests)
- ≡ • Optionally: Continuous deployment to production server

# Dashboards for Model Evaluation Results

2017-08-19-06-29-22-855-UTC

SUMMARY DEPLOY RETRAIN

PERFORMANCE MODEL VIS FEATURES

### Test Data Performance

threshold: 0.0584

0.288

0.925

performance

0.7936 auc

Precision-Recall

Recall: 0.69 Precision: 0.54

ROC

FPR: 0.252 TPR: 0.69

Confusion Matrix

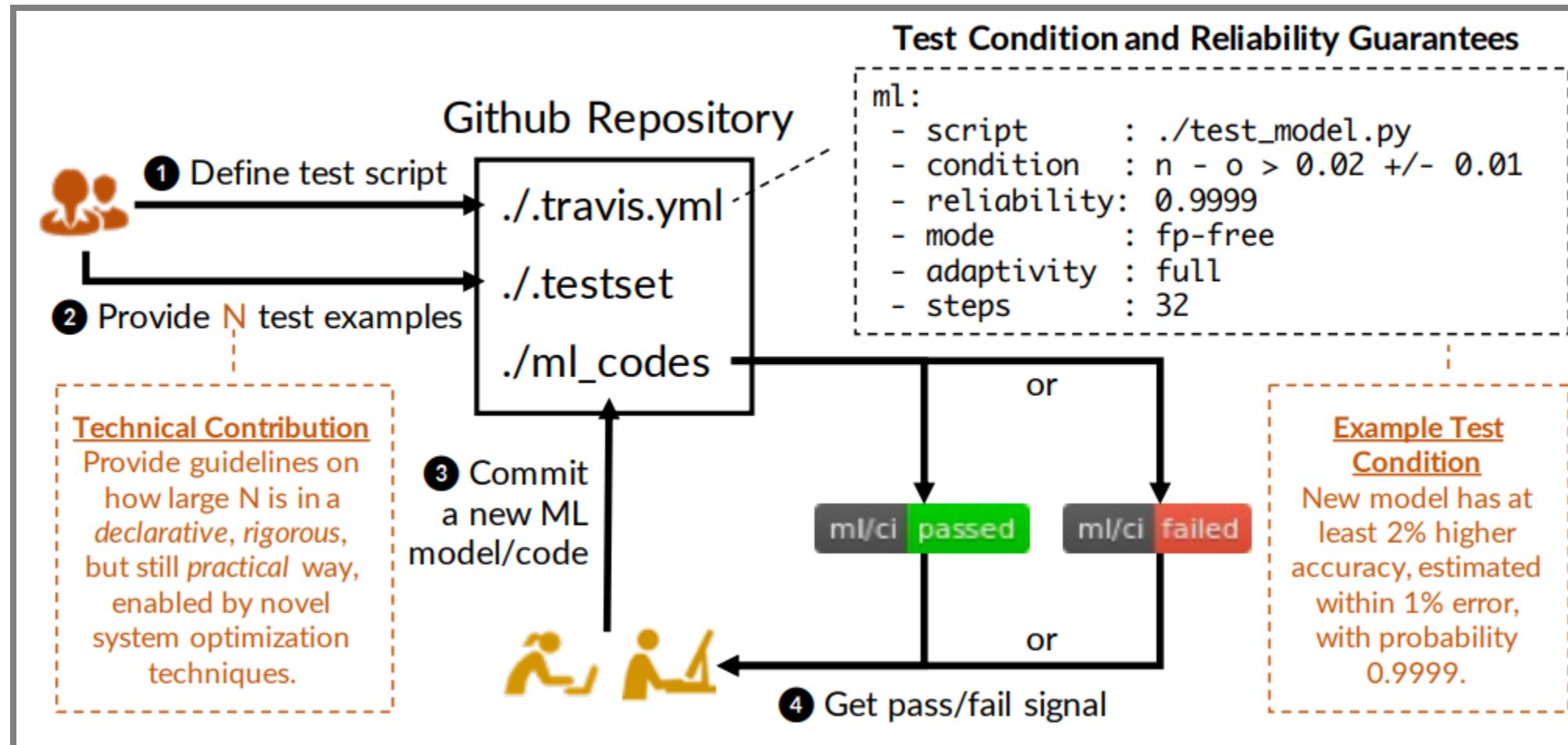
Positive label: true

Predicted		
Actual	YES	NO
YES	TP: 0.21 17604 Samples	FN: 0.093 7891 Samples
NO	FP: 0.18 15005 Samples	TN: 0.52 44549 Samples

### calibration reliability

The reliability diagram shows how reliable (or "well-calibrated") the model's probability estimates are when evaluated on the test data. For example, A well calibrated (binary) model should classify the samples such that among the samples to which it gives a probability close to 0.8 of belonging to the positive class, approximately 80% of those samples actually belong to

# Specialized CI Systems



# Dashboards for Comparing Models

# Summary

## Curating test data

- Analyzing specifications, capabilities
- Not all inputs are equal: Identify important inputs (inspiration from specification-based testing)
- Slice data for evaluation
- Identifying capabilities and generating relevant tests

## Automated random testing

- Feasible with invariants (e.g. metamorphic relations)
- Sometimes possible with simulation

Automate the test execution with continuous integration

# Further readings

- Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "[Semantically equivalent adversarial rules for debugging NLP models.](#)" In Proc. ACL, pp. 856-865. 2018.
- Barash, Guy, Eitan Farchi, Ilan Jayaraman, Orna Raz, Rachel Tzoref-Brill, and Marcel Zalmanovici. "[Bridging the gap between ML solutions and their business requirements using feature interactions.](#)" In Proc. FSE, pp. 1048-1058. 2019.
- Ashmore, Rob, Radu Calinescu, and Colin Paterson. "[Assuring the machine learning lifecycle: Desiderata, methods, and challenges.](#)" arXiv preprint arXiv:1905.04223. 2019.
- Christian Kaestner. [Rediscovering Unit Testing: Testing Capabilities of ML Models](#). Toward Data Science, 2021.
- D'Amour, Alexander, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen et al. "[Underspecification presents challenges for credibility in modern machine learning.](#)" arXiv preprint arXiv:2011.03395 (2020).
- Segura, Sergio, Gordon Fraser, Ana B. Sanchez, and Antonio Ruiz-Cortés. "[A survey on metamorphic testing.](#)" IEEE Transactions on software engineering 42, no. 9 (2016): 805-824.