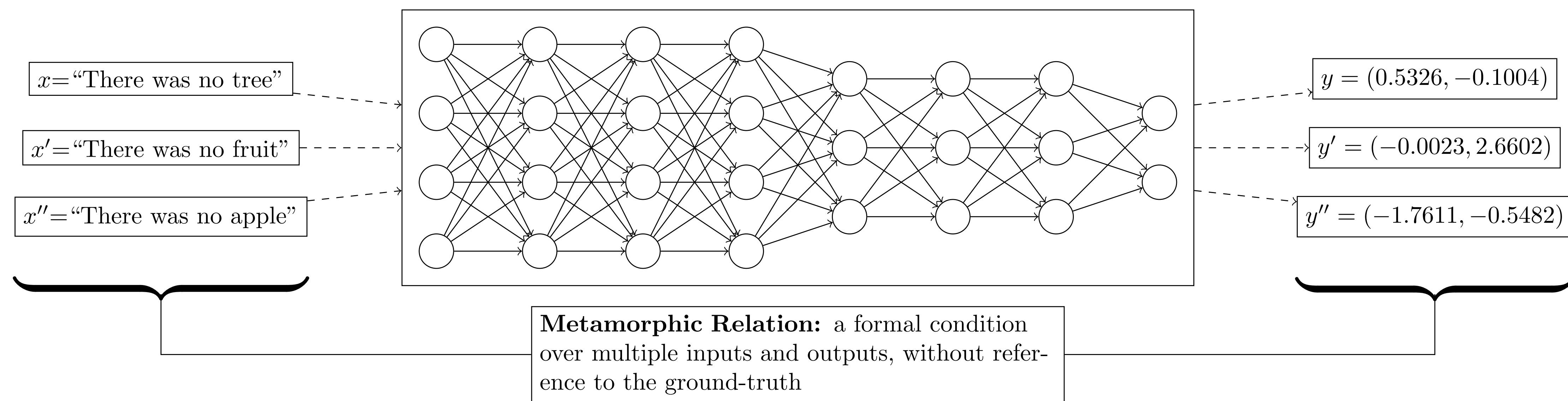


# Generate millions of test cases from few unlabelled data with metamorphic testing.

## Systematicity, Compositionality and Transitivity of Deep NLP Models: a Metamorphic Testing Perspective

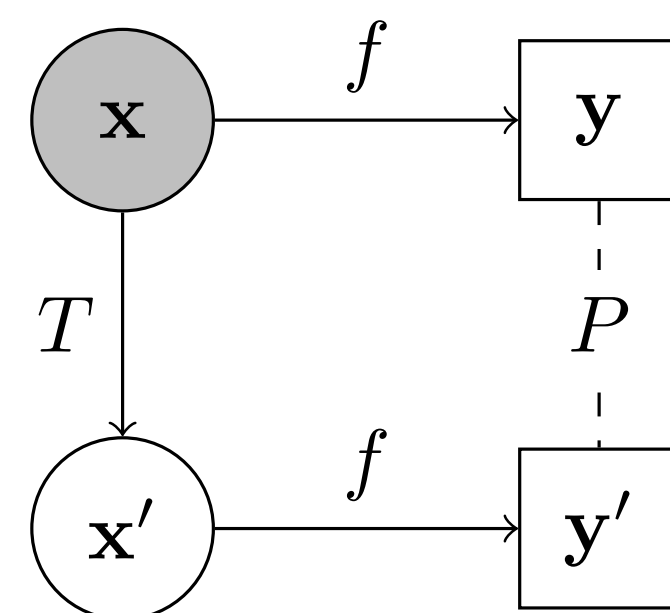
Edoardo Manino, Julia Rozanova, Danilo Carvalho, André Freitas, Lucas Cordeiro  
University of Manchester (UK), Idiap Research Institute (CH), EnnCore project



### Existing works: single-input metamorphic relations

Existing metamorphic testing for NLP:

- One input from test set
- Robustness-like relations
- $T \rightarrow$  typos, synonyms, etc.
- $P \rightarrow$  same output class



#### Single-input metamorphic test

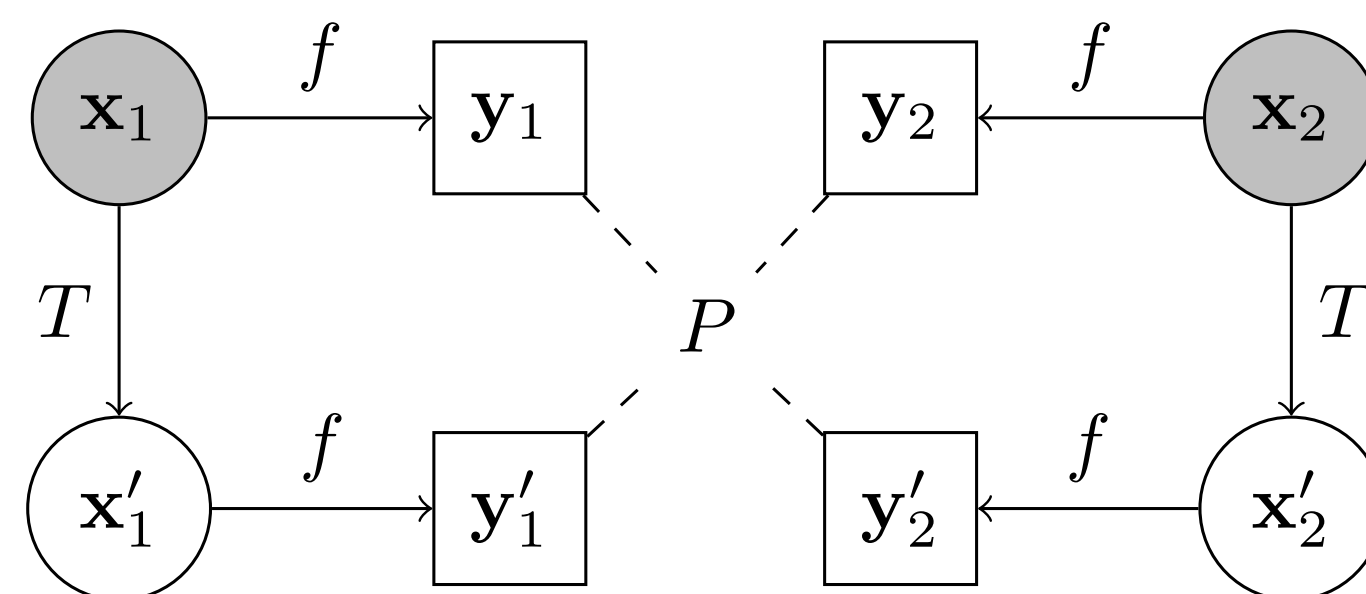
Input:  $x =$  The cat sat on the mat.  
 $x' =$  The pet stood onto the mat.  
 $T:$  replace any word of the input with a synonym.  
 $P:$   $y = f(x) \wedge \exists i \forall j \neq i (y_i > y_j) \wedge (y'_i > y'_j)$

Is the output class preserved after replacing some input words with synonyms?

### Contribution 1: pairwise systematicity relations

Test internal consistency of model:

- Two inputs from test set
- Read their output relation
- Is it preserved after applying  $T$ ?
- 112M+ tests from 11K+ data
- RoBERTa sentiment: 5-10% errors



#### Pairwise systematicity metamorphic

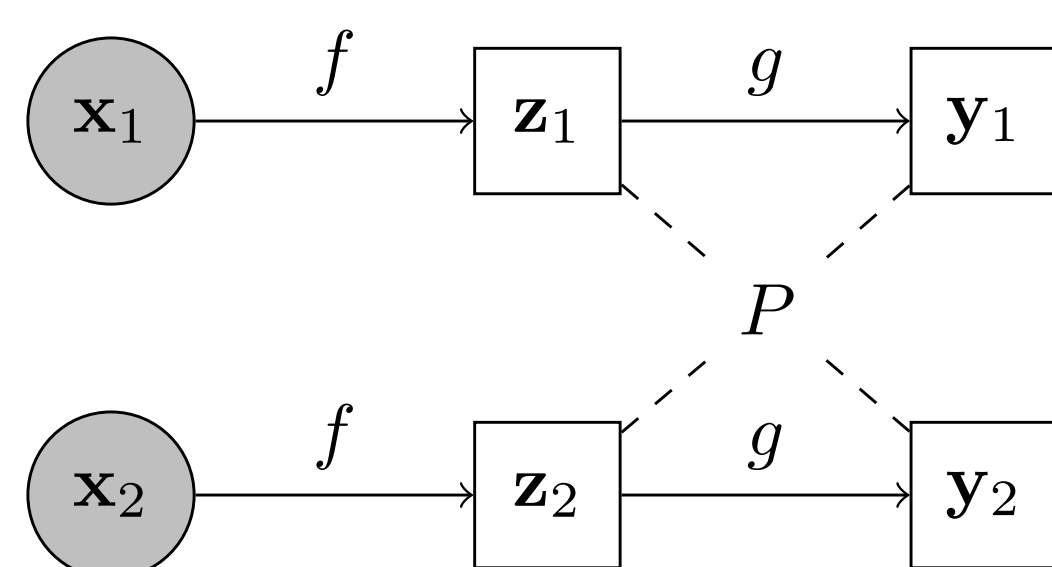
Input:  $x_1 =$  Light, cute and forgettable.  
 $x_2 =$  A masterpiece four years in the making.  
 $x'_1 =$  Thank you. Light, cute and forgettable.  
 $x'_2 =$  Thank you. A masterpiece four years in the making.  
 $T:$  concatenate the text Thank you. at the beginning of the input.  
 $P:$   $s_{pos}(f(x_1)) > s_{pos}(f(x_2)) \iff s_{pos}(f(x'_1)) > s_{pos}(f(x'_2))$

Is the polarity between the two sentences preserved after concatenation of the fragment?

### Contribution 2: pairwise compositionality relations

Metamorphic version of probing:

- Two inputs from test set
- Probe hidden reps. after  $f$
- Does it correlate with output?
- 9M+ tests from fewer than 1K data
- RoBERTa entailment: 25-70% errors



#### Pairwise compositionality metamorphic

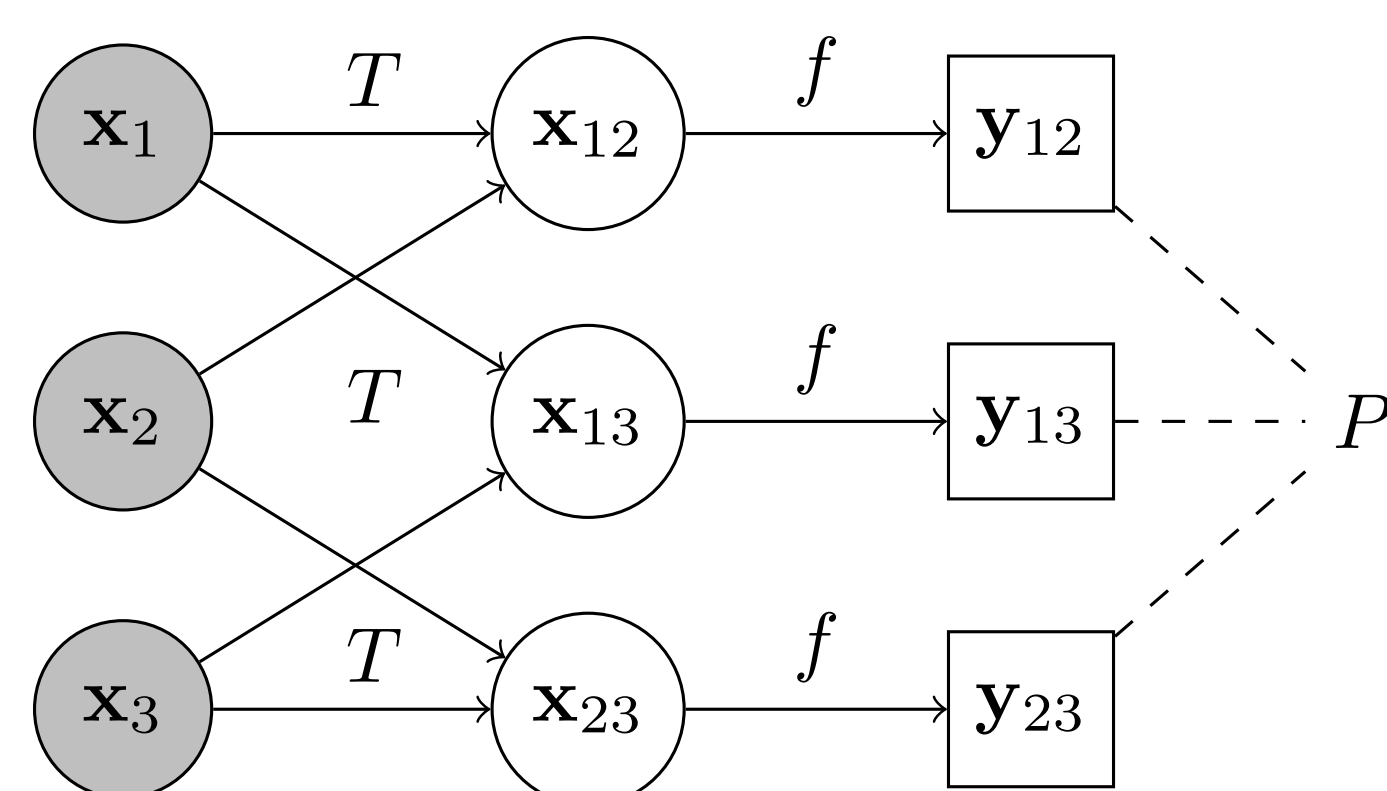
Input:  $x_1 =$  There was no tree. There was no cherry tree.  
 $x_2 =$  There was no fruit. There was no apple.  
Hidden:  $f(x_1) =$  contextual embeddings of the tokens ( tree. cherry tree. )  
 $f(x_2) =$  contextual embeddings of the tokens ( fruit. apple. )  
 $P:$   $s_{hyp}(f(x_1)) > s_{hyp}(f(x_2)) \iff s_{ent}(g(f(x_1))) > s_{ent}(g(f(x_2)))$

Does the polarity between the two embeddings correspond to the polarity of the output?

### Contribution 3: three-way transitivity relations

Are mistakes transitive too?

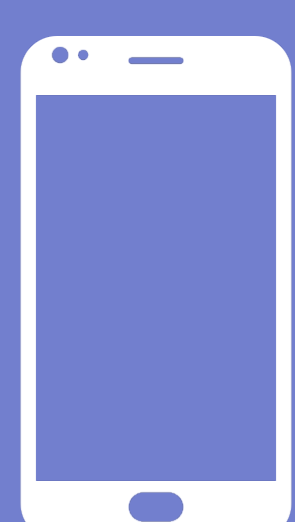
- Three inputs from test set
- If two pairs are predicted true...
- ...the third must be true too!
- Cubic number of test cases
- RoBERTa lex. rel.: 60-80% errors



#### Three-way transitivity metamorphic

Input:  $x_1, x_2, x_3 =$  arrangement symmetrical together  
 $x_{12} =$  ( arrangement symmetrical )  
 $x_{23} =$  ( symmetrical together )  
 $x_{13} =$  ( arrangement together )  
 $T:$  choose two words from the source triplet  $x_1, x_2, x_3$   
 $P_{syn}: v_{syn}(f(x_{12})) \wedge v_{syn}(f(x_{23})) \implies v_{syn}(f(x_{13}))$   
 $P_{hyp}: v_{hyp}(f(x_{12})) \wedge v_{hyp}(f(x_{23})) \implies v_{hyp}(f(x_{13}))$

When the model classifies two input pairs as positive, does it also classify the third as positive?



Scan QR code to get the full paper

