# MERGE: MINIMAL EXPRESSION REPLACEMENT GENERALIZATION TEST

MĂDĂLINA ZGREABĂN          TEJASWINI DEOSKAR          LASHA ABZIANIDZE

Utrecht University

## NLI BENCHMARKS

- Disturb lexical overlap heuristic of premise and hypothesis (*PH*);
- Have low lexical diversity;
- Costly, if formed manually;
- Syntax non-preserving;
- Unfair, if the data is not similar enough to the training data.
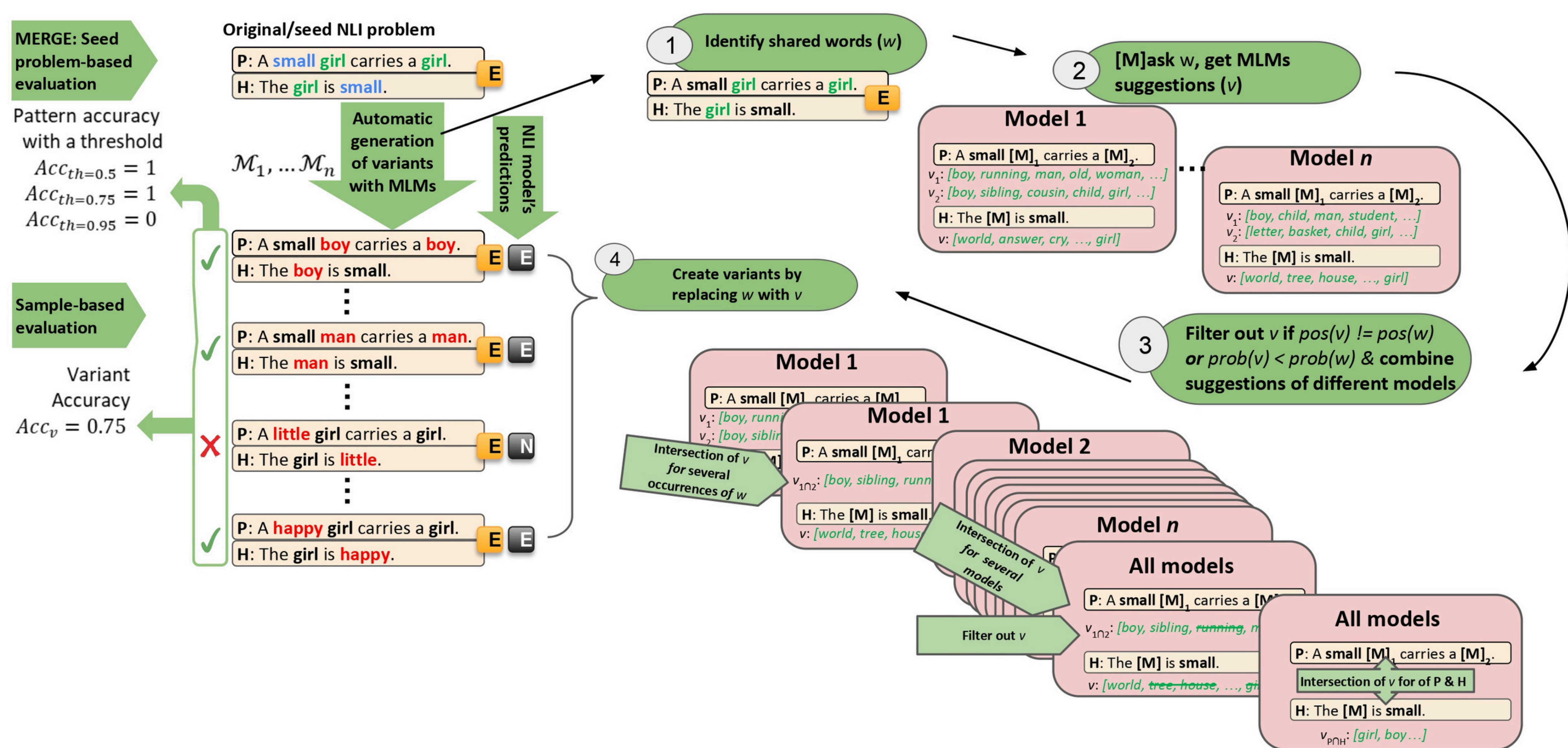
## MERGE

- Minimally alters already existing NLI datasets;
- Preserves underlying logical reasoning;
- Does not require human validation by strict minimal changing criteria;
- Preserves lexical overlap;
- Can add more lexical diversity by adding suggestions from other models;
- Automatic;
- Syntax preserving.

## RESEARCH QUESTIONS
### ARE LANGUAGE MODELS ROBUST AGAINST MINIMAL VARIANTS OF NLI PROBLEMS?
### DO THE LIKELIHOOD, POS TAG, PLAUSIBILITY, OR MASKED MODELS MATTER?
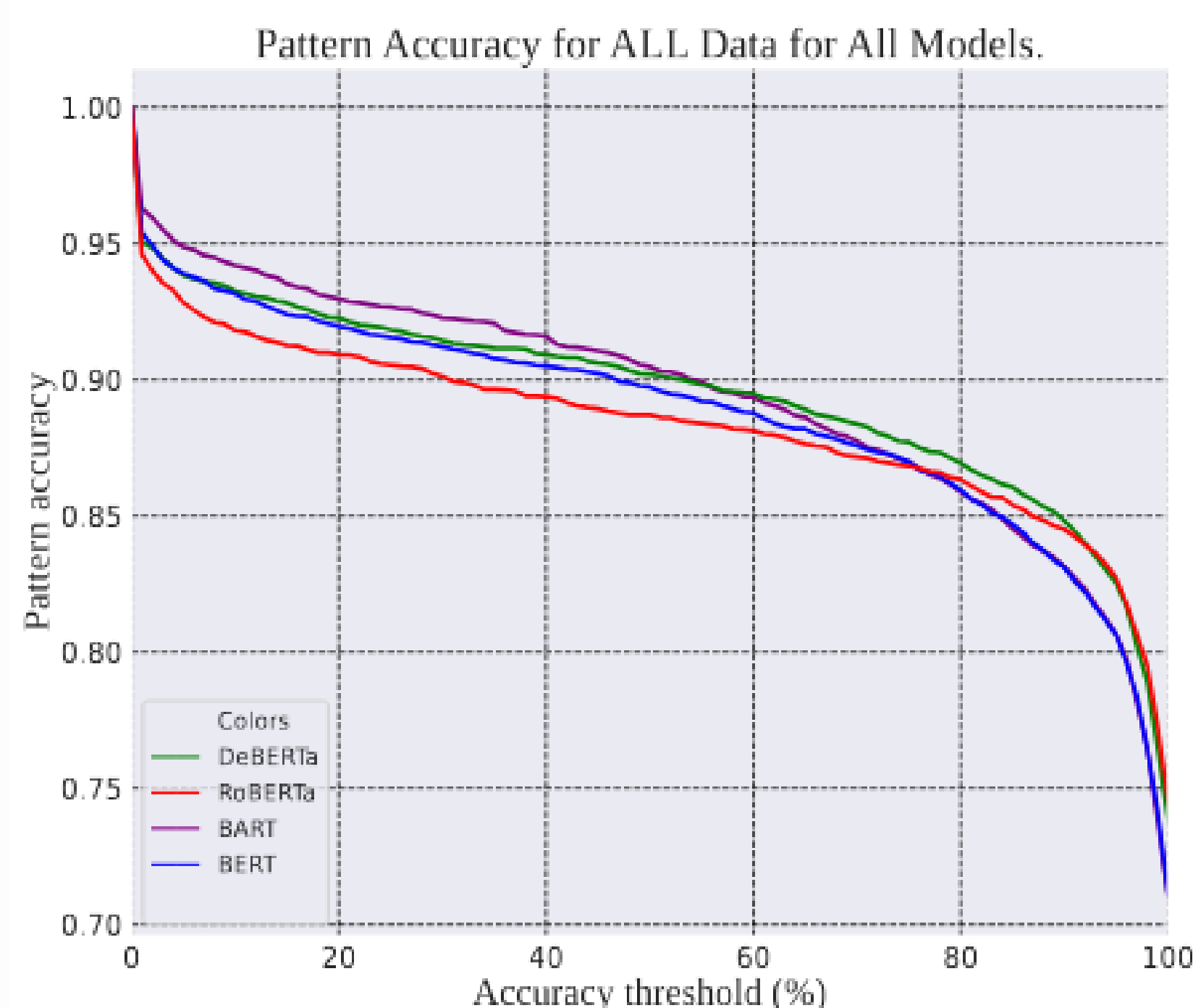
## FRAMEWORK



## METHODOLOGY

- 200 suggestions (*v*) with bert-base-cased and roberta-base;
- Suggestions tagged with en_core_web_sm;
- Exclude punctuation signs, derviational morephems, different POS(*v*), and lower probability(*v*).
- Required variants ==20;
- 10 random mini-datasets with 20 variants per problem (ALL_Var).
- Evaluate BERT, BART, DeBERTa, RoBERTa.

## THE LIKELIHOOD, POS TAG, PLUASIBILITY & MLMS DO MATTER

## LMS ARE NOT ROBUST AGAINTS MINIMAL VARIANTS



Pattern Accuracy for ALL Data for All Models.

| Word | Seed | Average | N(%) | C(%) | E(%) | Subs |
|---|---|---|---|---|---|---|
| $N_{Var}$ | 3704 | 144.2 | 12.5 | 22.6 | 46.1 | 74080 |
| $V_{Var}$ | 1129 | 112 | 28.1 | 16.6 | 55.2 | 22580 |
| $Adj_{Var}$ | 280 | 79.9 | 32.5 | 22.5 | 44.8 | 5620 |
| $ALL_{Var}$ | 4468 | 152.8 | 30.7 | 21.4 | 47.7 | 102280 |

| Model | Training | $SNLI_{test}$ | $ALL_{Seed}$ | $ALL_{Variants}$ |
|---|---|---|---|---|
| BERT | S | 90.48 | 90.24 | 88.72 |
| RoBERTa | S | 90.06 | 89.86 | 88.50 |
| DeBERTa | S | 91.70 | 91.38 | 89.41 |
| BART | S, M, F, A | 92.03 | 91.85 | 89.11 |

| Model | Training | $All_{BERT}$ | $All_{RoBERTa}$ | $All_{Both}$ |
|---|---|---|---|---|
| BERT | S | 88.79 | 88.55 | 88.84 |
| RoBERTa | S | 88.58 | 88.33 | 88.56 |