# Part-of-Speech Annotation in 4 hours: How to spend your time
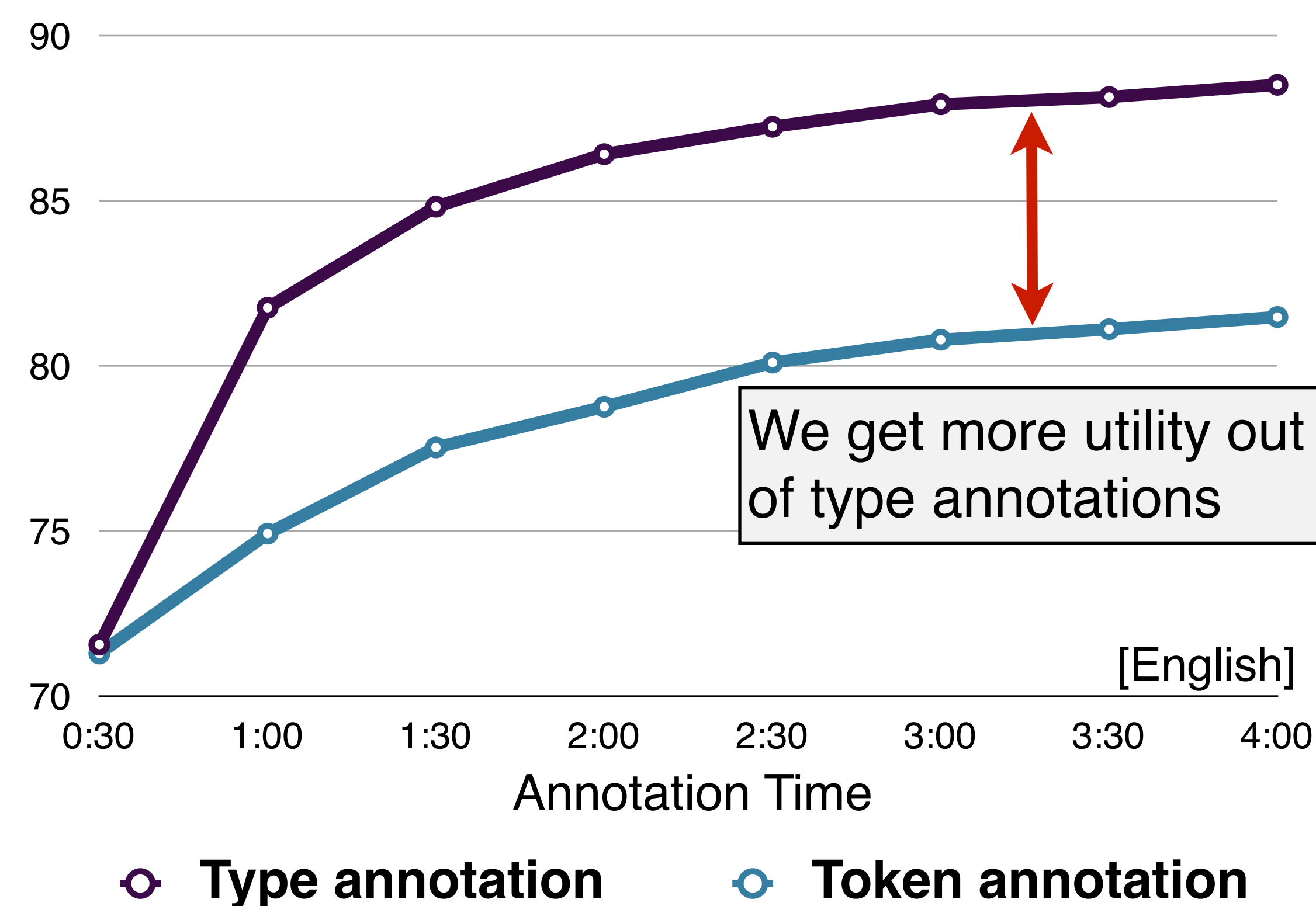
## Dan Garrette, Jason Mielens, and Jason Baldridge

The University of Texas at Austin

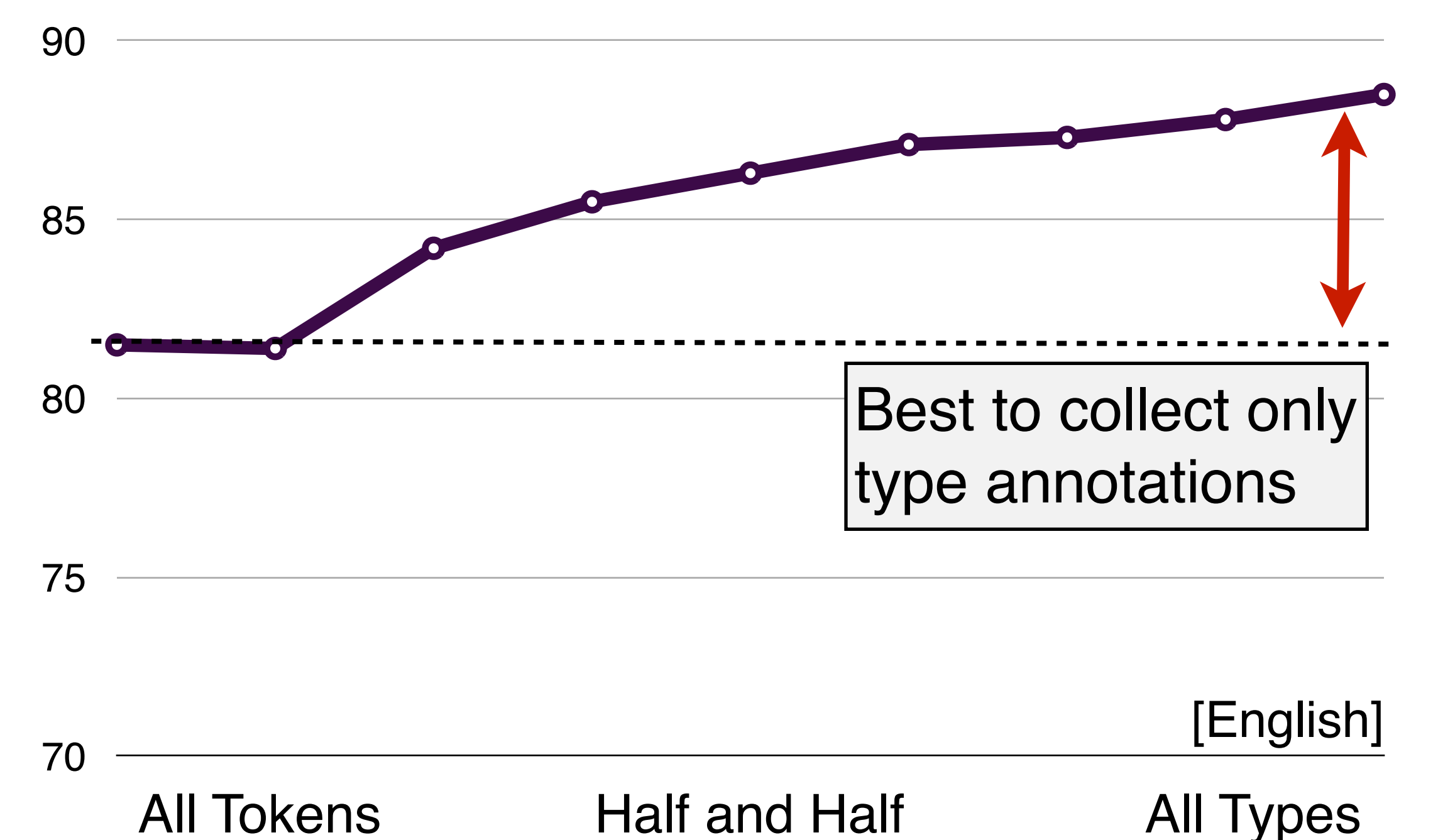For low-resource languages, we cannot expect much training data.

- Can we train with less data?

- What can we collect cheaply?

- Where to focus collection efforts?

We conducted experiments collecting data from linguists in a fixed time.
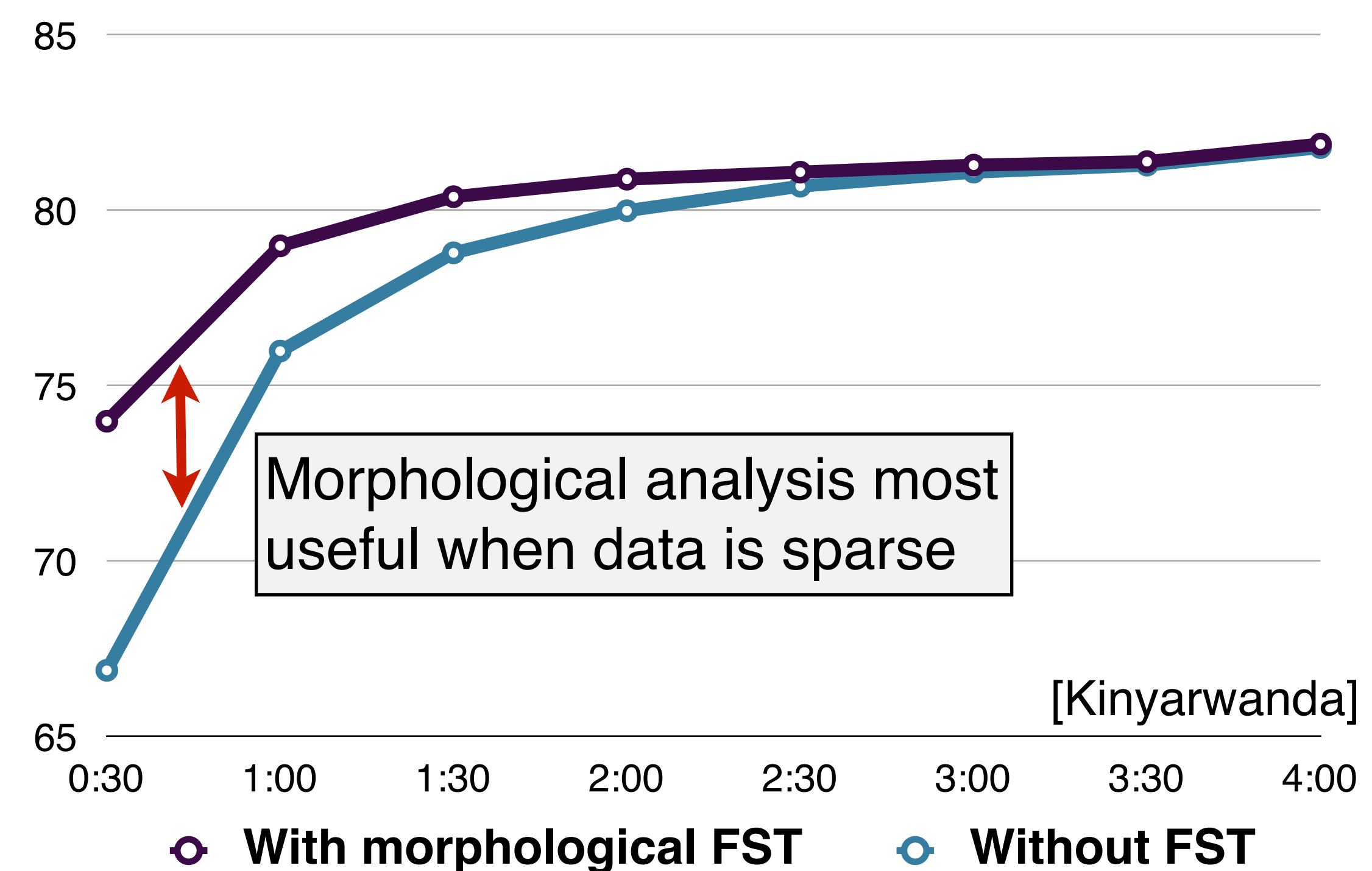
Our semi-supervised approach:

- Expand the initial annotations

- Remove expansion noise by finding a minimal model

- Train a MaxEnt Markov Model

Achieve **90%** accuracy after **4 hours**

[Garrette & Baldridge, NAACL 2013]

## Annotating: Types > Tokens



We get more utility out of type annotations

[English]

Type annotation    Token annotation

## Mixing Type and Token Annotations



Best to collect only type annotations

[English]

## Morphological analysis with an FST helps for Kinyarwanda



Morphological analysis most useful when data is sparse

[Kinyarwanda]

With morphological FST    Without FST

## Annotator Experience



Experience advantage is greater with types

[English]

Experienced Annotator    Novice Annotator