| Methods | Distillation | Latent Variables | Latent Alignments | Glancing Targets |
|---|---|---|---|---|
| What it can do? | simplifying the training data | model any types of dependency in theory | handling token shifts in the output space | ease the difficulty of learning hard examples |
| What it cannot? | uncertainty exists in the teacher model | constrained by the modeling power of the used latent variables | unable to model non-monotonic dependency, e.g. reordering | training / testing phase mismatch |
| Potential issues | sub-optimal due to the teacher's capacity | difficult to train; posterior collapse | decoder inputs must be longer than targets | difficult to find the optimal masking ratio |