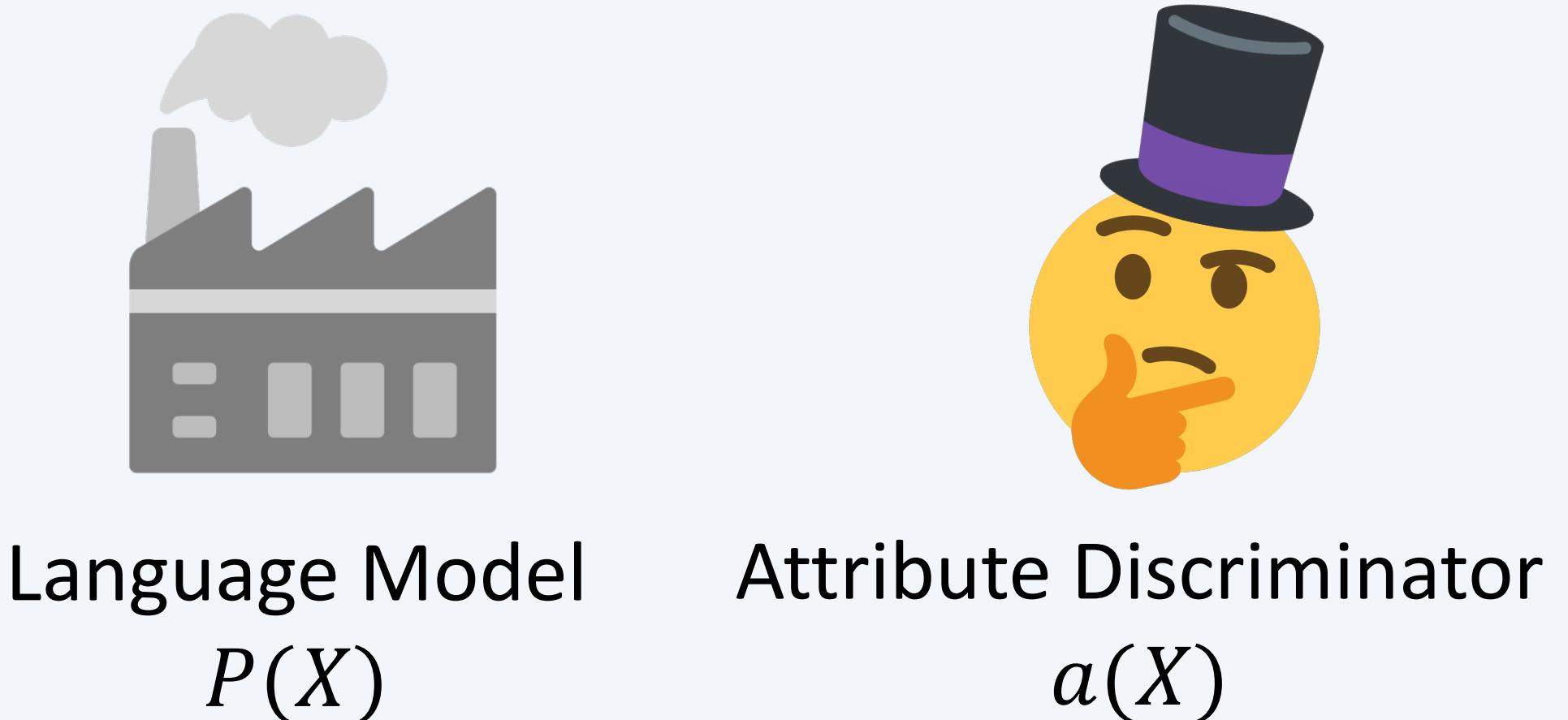


FUDGE: Controlled Text Generation With Future Discriminators

Kevin Yang and Dan Klein
UC Berkeley

Controlled Text Generation



For a desired attribute a (e.g. formal style), controlled text generation samples from the distribution $P(X|a(X) = \text{True})$, or henceforth $P(X|a)$.

Evaluation Tasks

Complete a Poem

And even thence thou wilt
be stol'n, I fear,

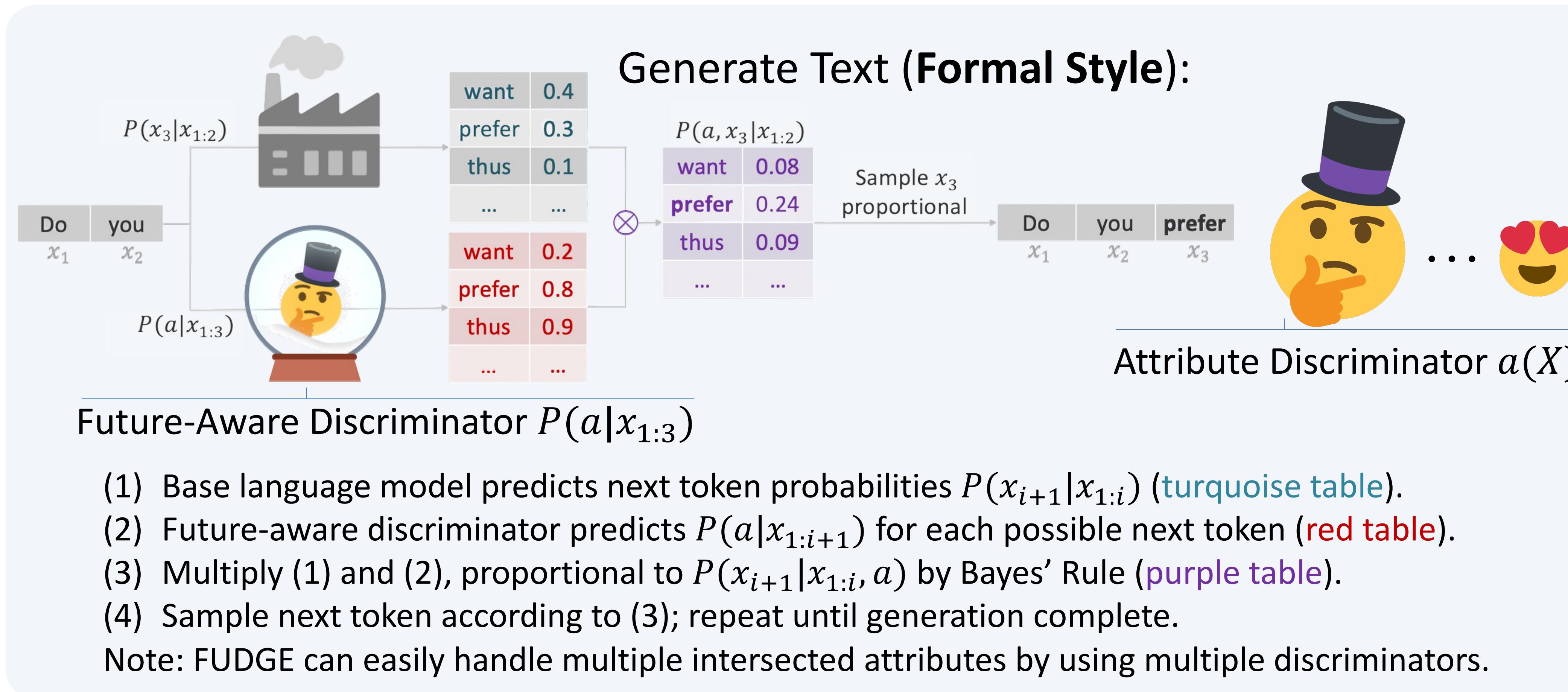
Generate Topical Document (Politics):

The issue focused on

Translate to English (Formal Style):

[la información que que
tenemos todo es propaganda
entonces es portante ver
otros versiones de lo que
está pasando en el mundo]

Guiding Generation Via Bayesian Factorization



Uncontrolled Examples

Language Model

And even thence thou wilt
be stol'n, I fear,
and thou art a good friend
of mine. The King's
Guard...

Language Model

The issue focused on
the role of a small company
called Korn Ferry, which
operates two vessels, the
Korn and...

Language Model

the information that, that
we have is all propaganda,
then is important to see
other versions of what's
happening...

FUDGE Examples

FUDGE

And even thence thou wilt
be stol'n, I fear,
for this shall be the end.
That's pretty clear.

FUDGE

The issue focused on
whether the two
institutions were
operating within the
bounds set by the
constitution...

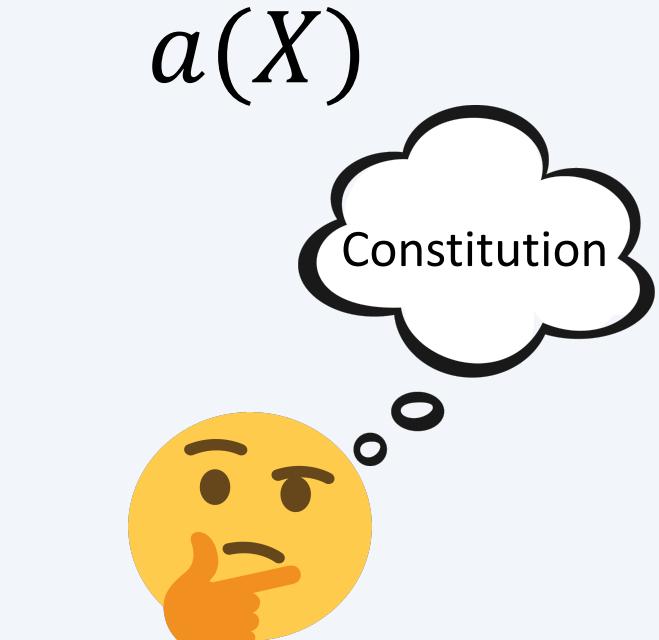
FUDGE

The information that we
have is all propaganda,
so, it's important to see
other versions of what's
happening...

Future Discriminators

Generate Text (Use "Constitution"):

Non-Future-Aware
Discriminator



Grep for "constitution"

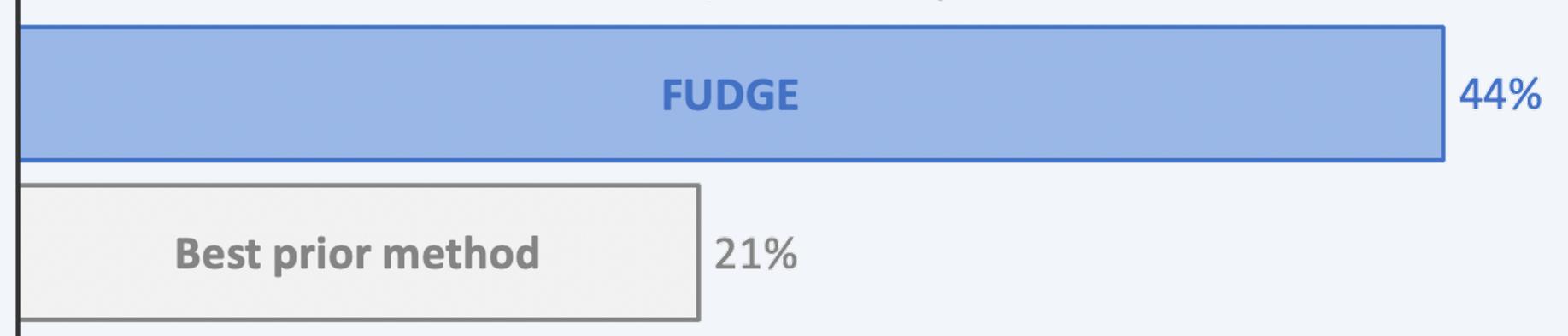
Future-Aware
Discriminator



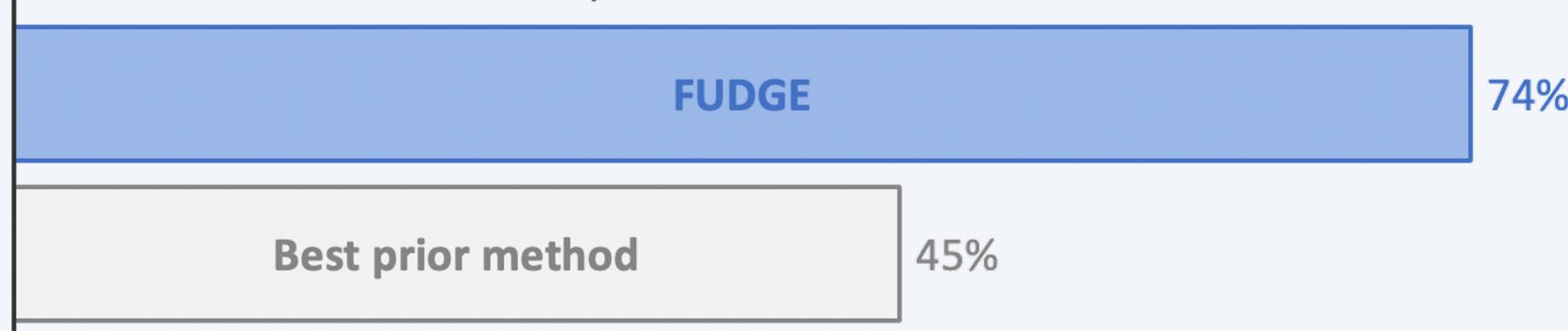
Does this continuation make
"constitution" more likely
later?

Results

% successful (iambic + rhyme + end sentence)



% on topic, human evaluations



Predicted probability of formality, avg. over test

