

limitations of autoregressive models and their alternatives

Chu-Cheng Lin[#], Aaron Jaech[‡], Xin Li[#], Matt Gormley[‡], Jason Eisner[#]

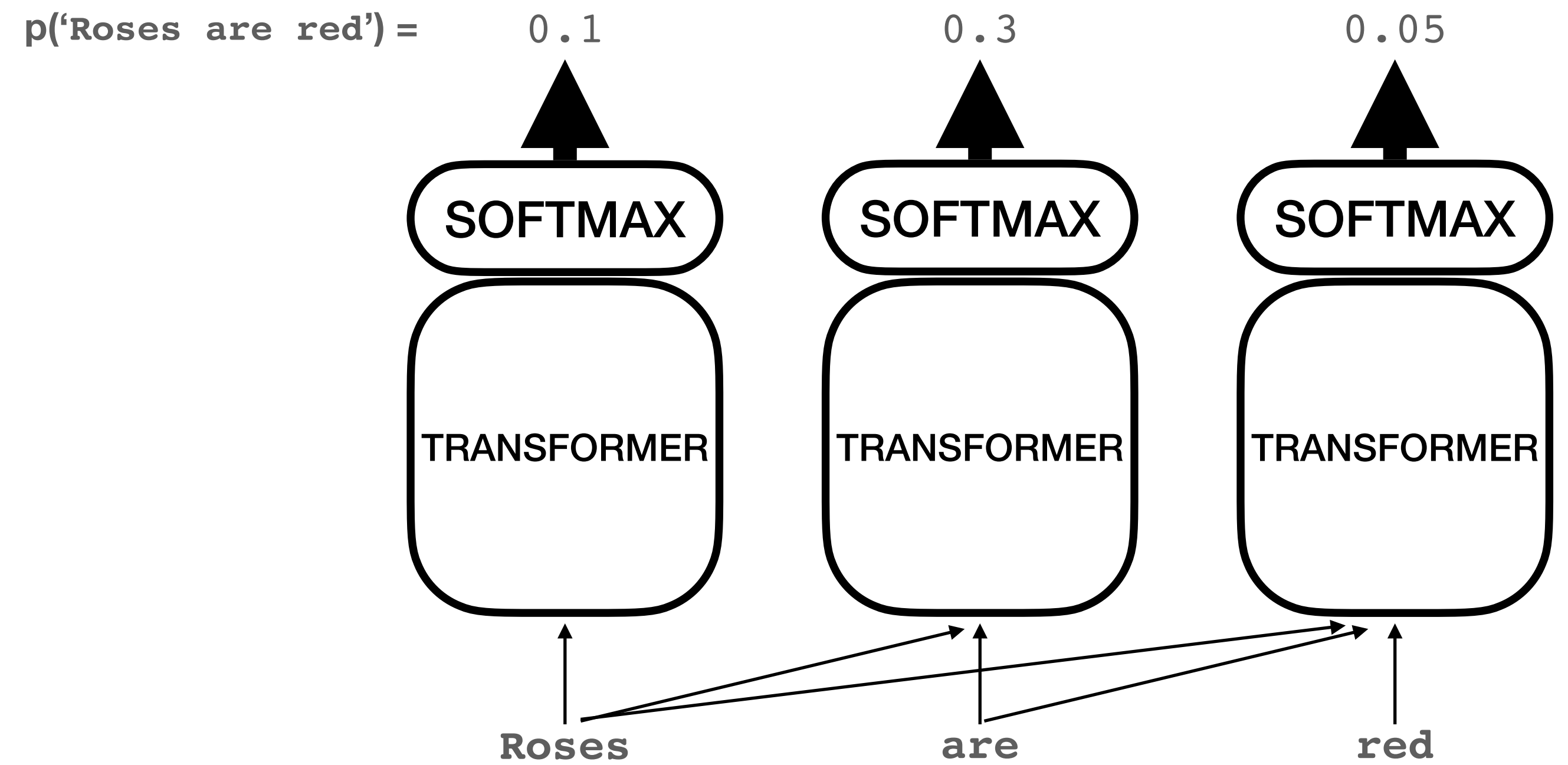
[#]Johns Hopkins University

[‡]Facebook AI

[‡]Carnegie Mellon University

This talk

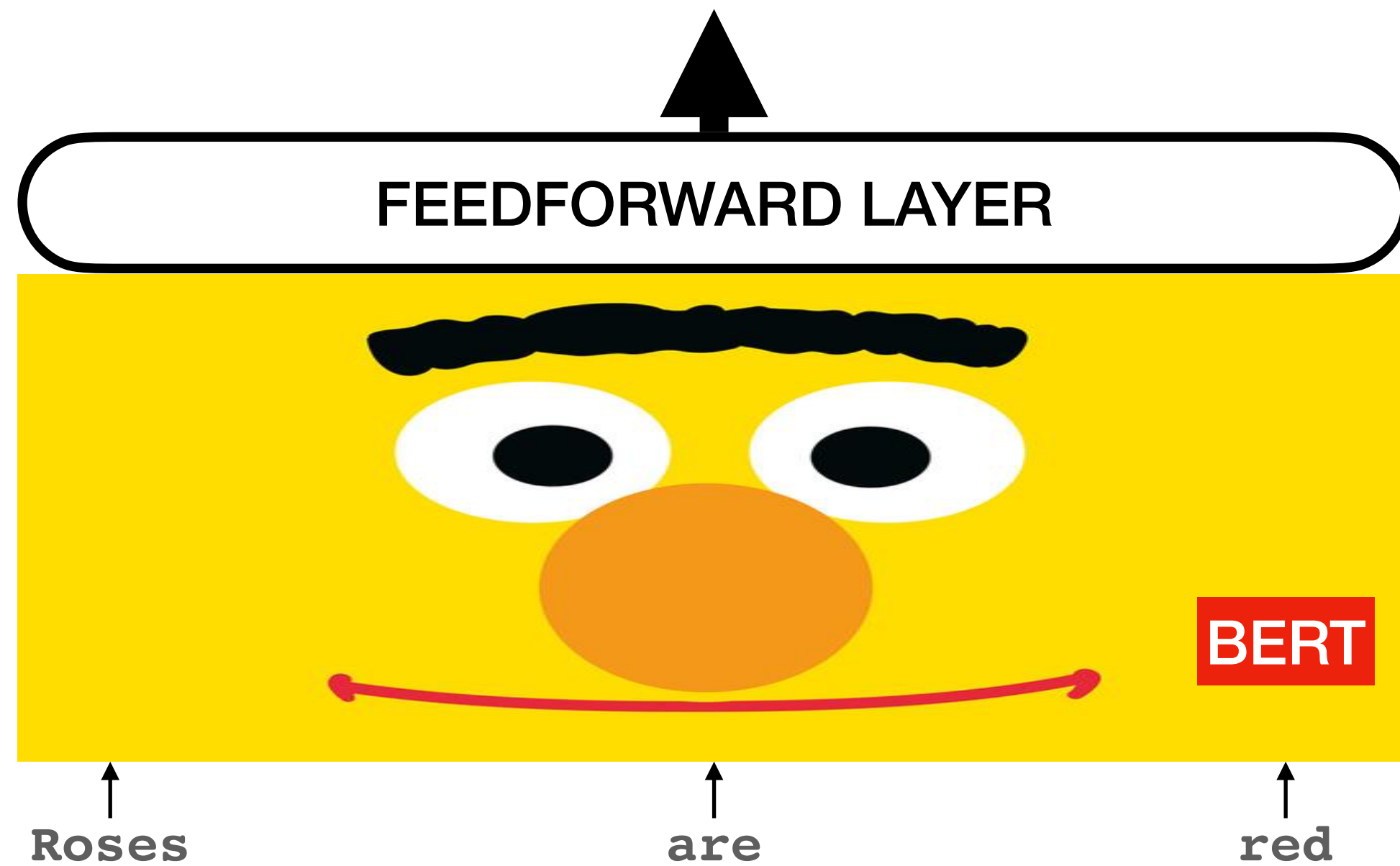
This talk



autoregressive models

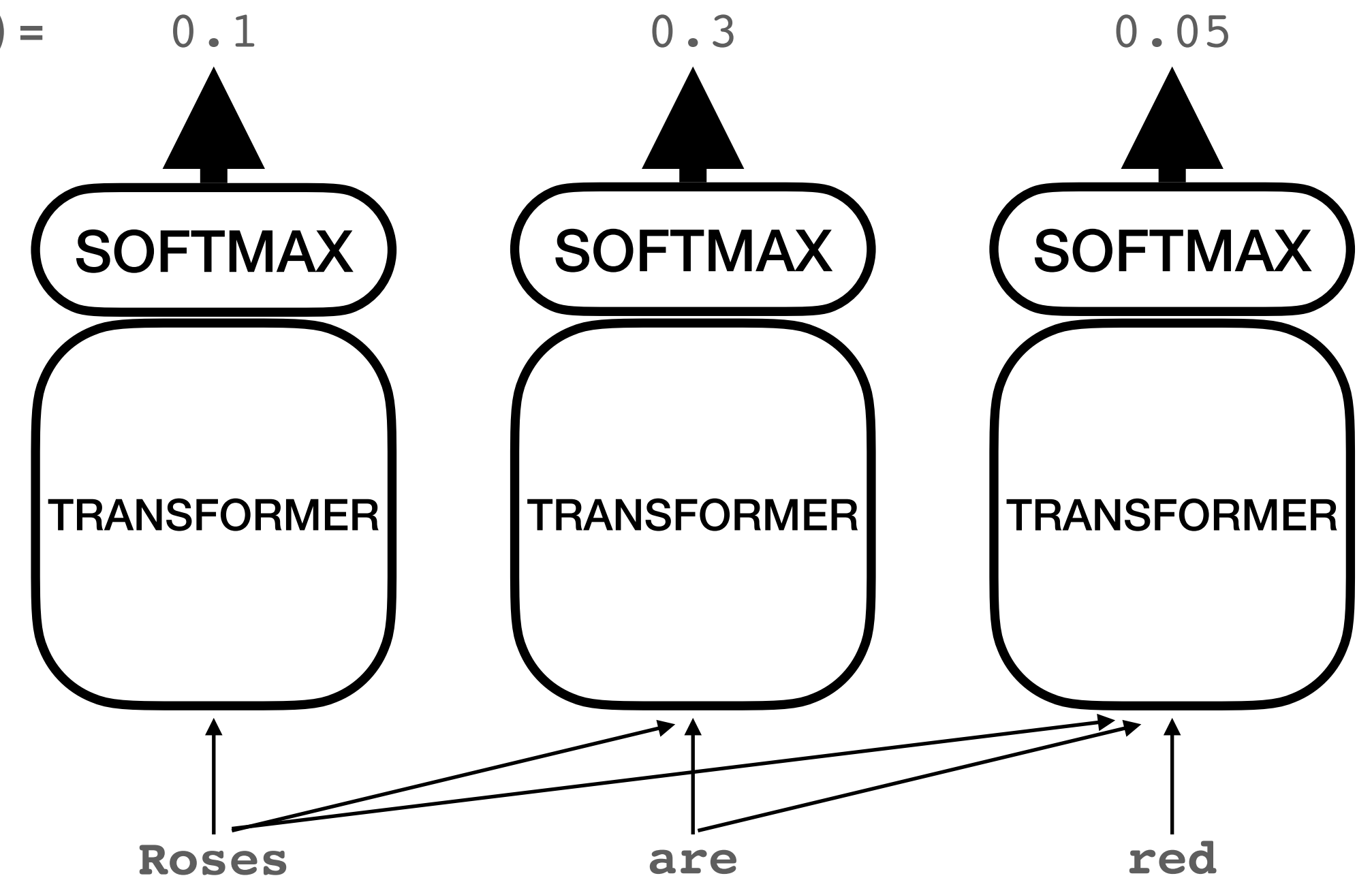
This talk

goodness('Roses are red') = 100



energy-based models (EBMs)

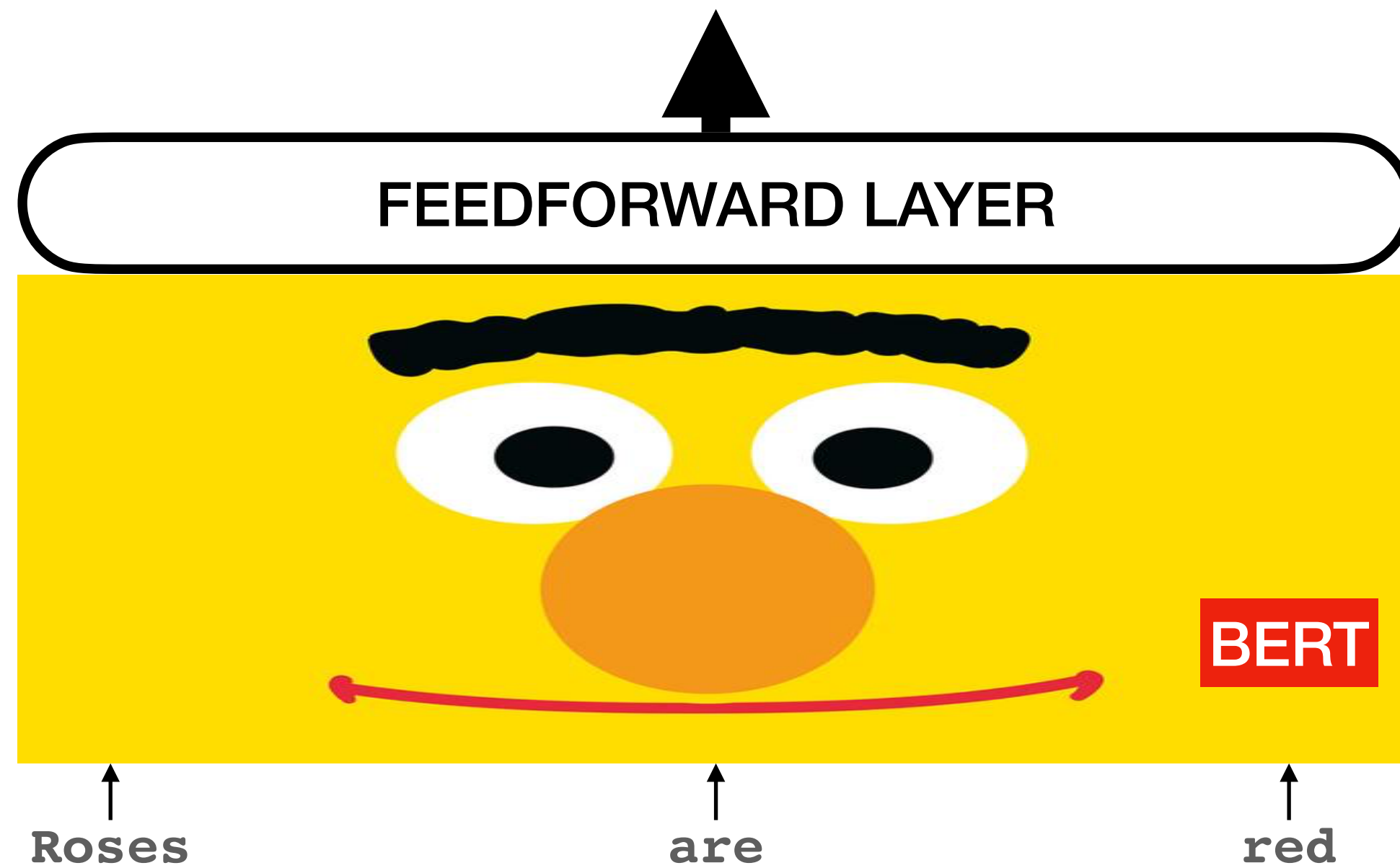
$p(\text{'Roses are red'}) =$



autoregressive models

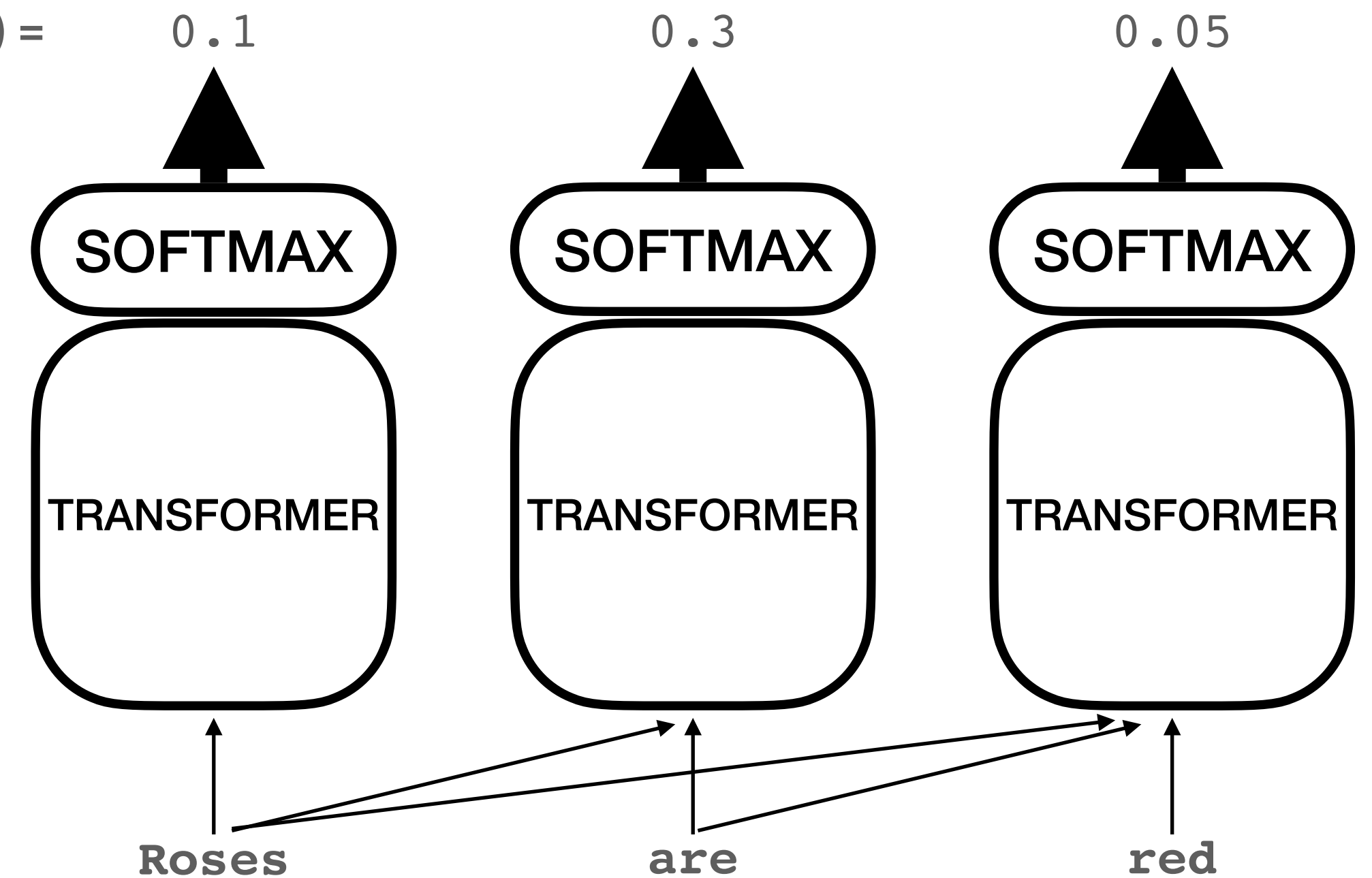
This talk

goodness('Roses are red') = 100



energy-based models (EBMs)

$p(\text{'Roses are red'}) =$

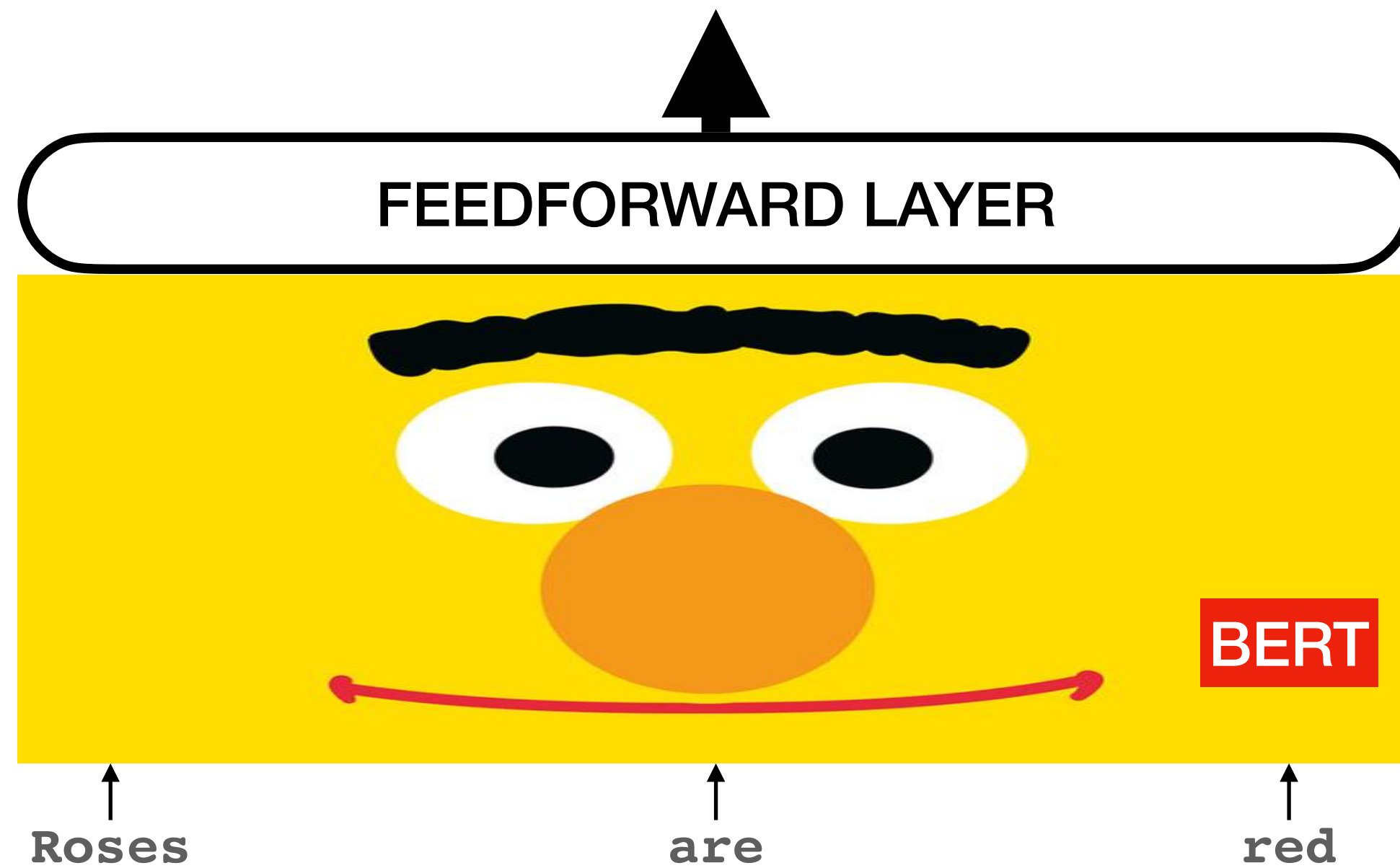


autoregressive models

**?
=**

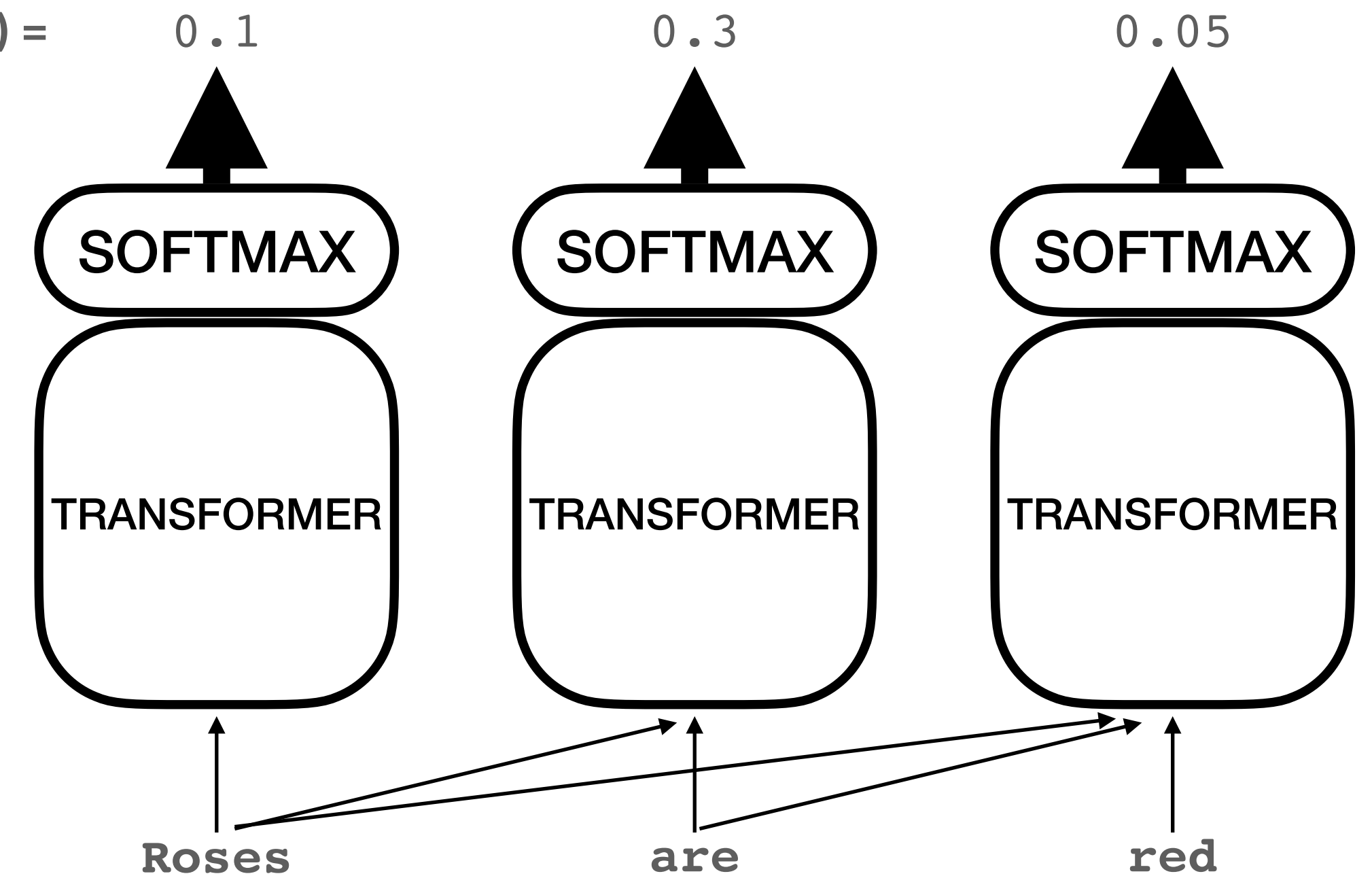
This talk

goodness('Roses are red') = 100

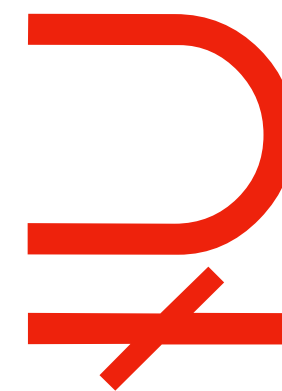


energy-based models (EBMs)

$p(\text{'Roses are red'}) =$



autoregressive models



Outline

- Autoregressive models are not as expressive as other model families, energy-based models in particular.
 - And having more parameters helps little!
- Model families that are more expressive than autoregressive models made their own trade-offs

Outline

- **Autoregressive models are not as expressive as other model families, energy-based models in particular.**
 - **And having more parameters helps little!**
- Model families that are more expressive than autoregressive models made their own trade-offs

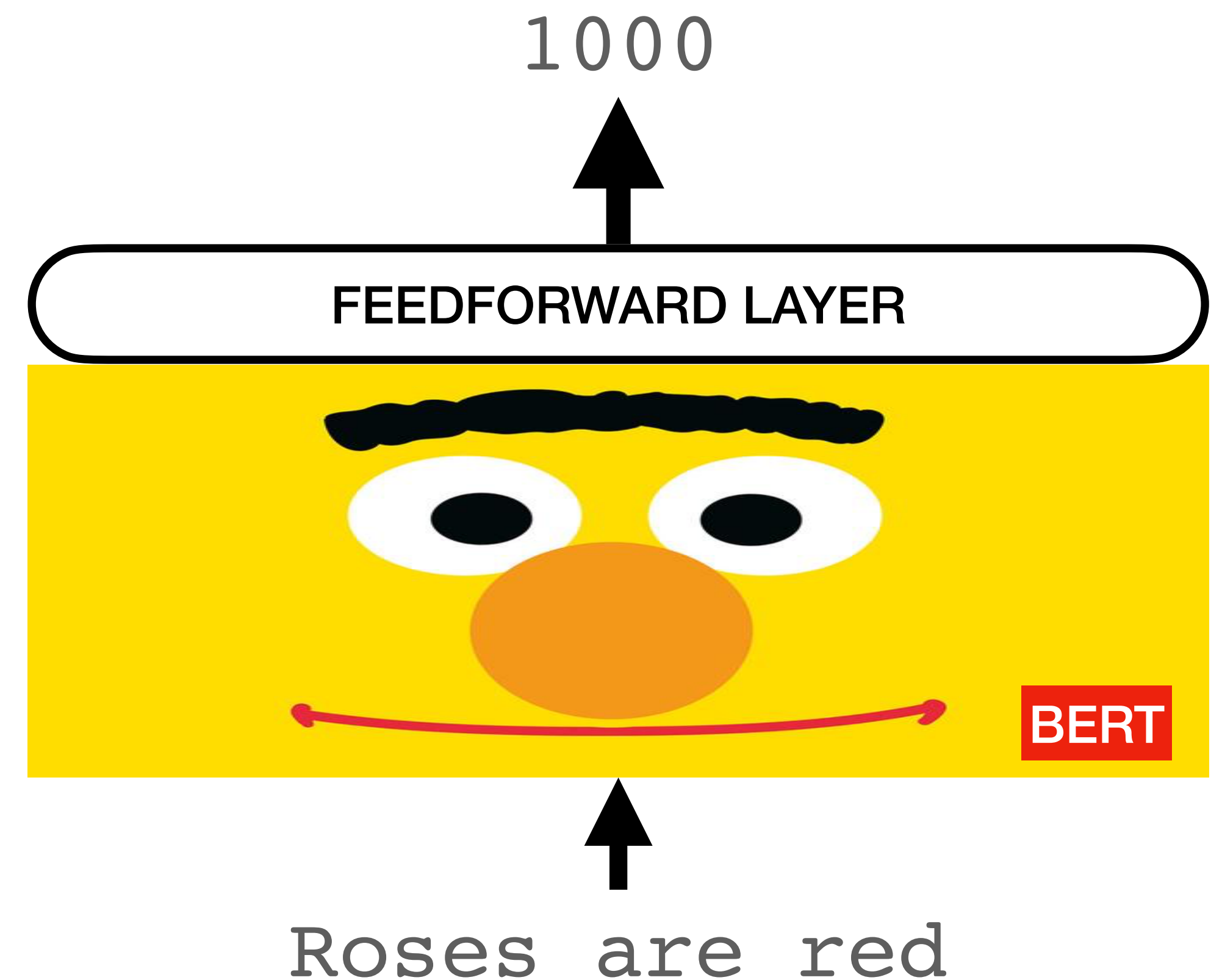
Many NLP tasks are about scoring strings

- Language modeling
 - **Good:**
Roses are red
 - **Maybe:**
Roses are nosy
 - **Bad:**
Roses queen sierra
- Machine translation
 - **Good:**
Roses are red -> Las
rosas son **rojas**
 - **Bad:**
Roses are red -> Las
rosas son **rojos**

Many NLP tasks are about scoring strings

we want to measure their goodness quantitatively with an NN

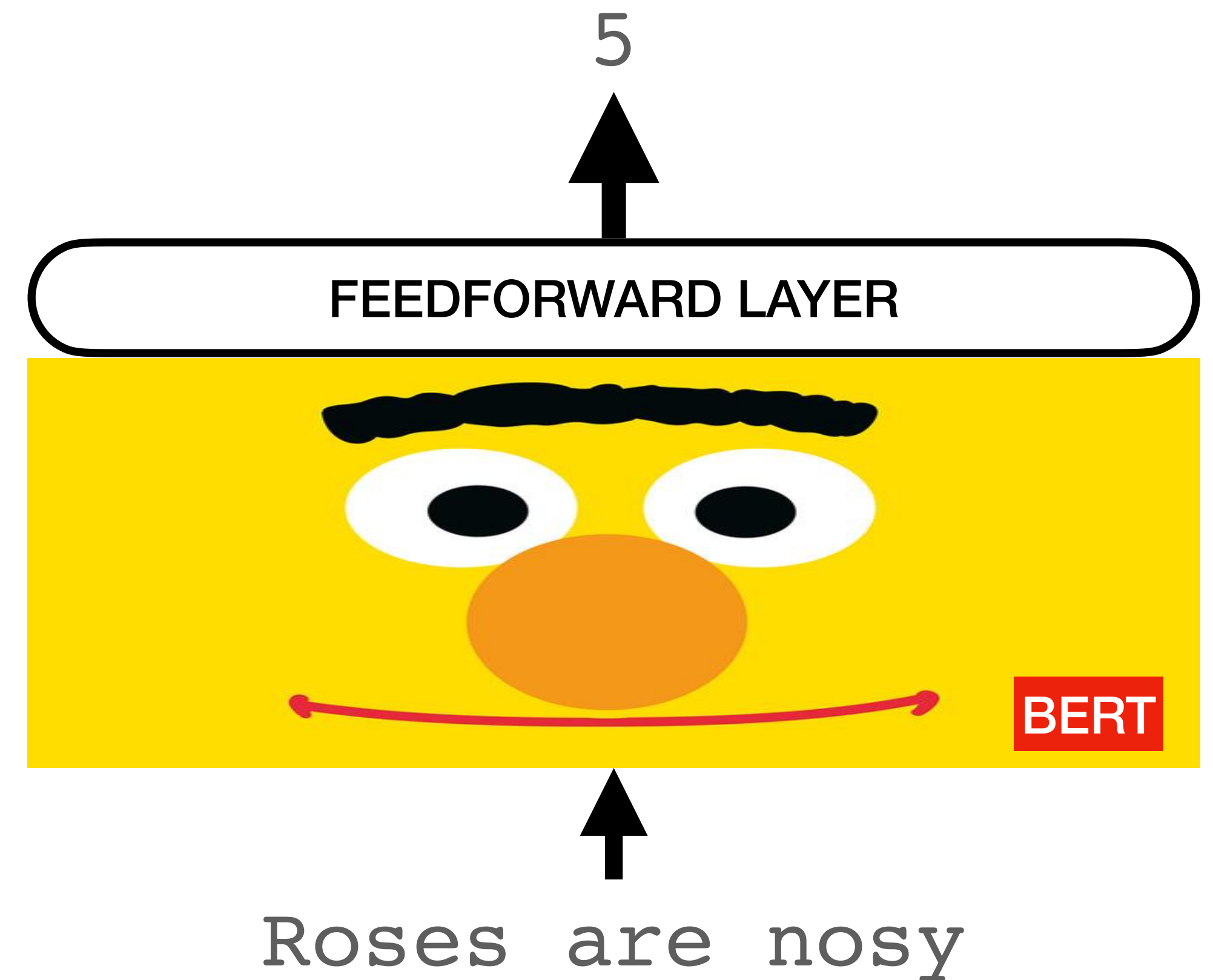
- **Good:**
Roses are red
- **Maybe:**
Roses are nosy
- **Bad:**
Roses queen sierra



Many NLP tasks are about scoring strings

we want to measure their goodness quantitatively with an NN

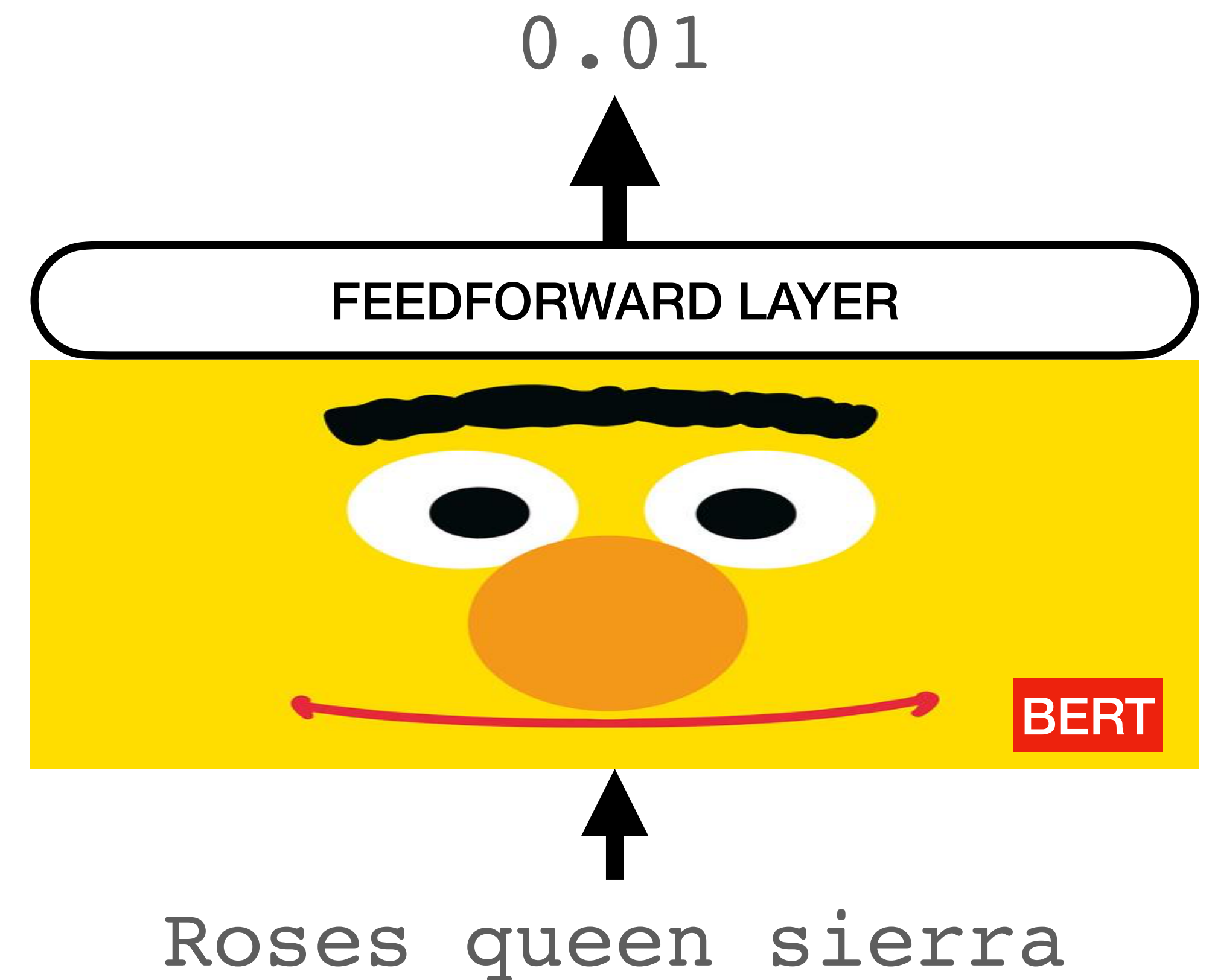
- **Good:**
Roses are red
- **Maybe:**
Roses are nosy
- **Bad:**
Roses queen sierra



Many NLP tasks are about scoring strings

we want to measure their goodness quantitatively with an NN

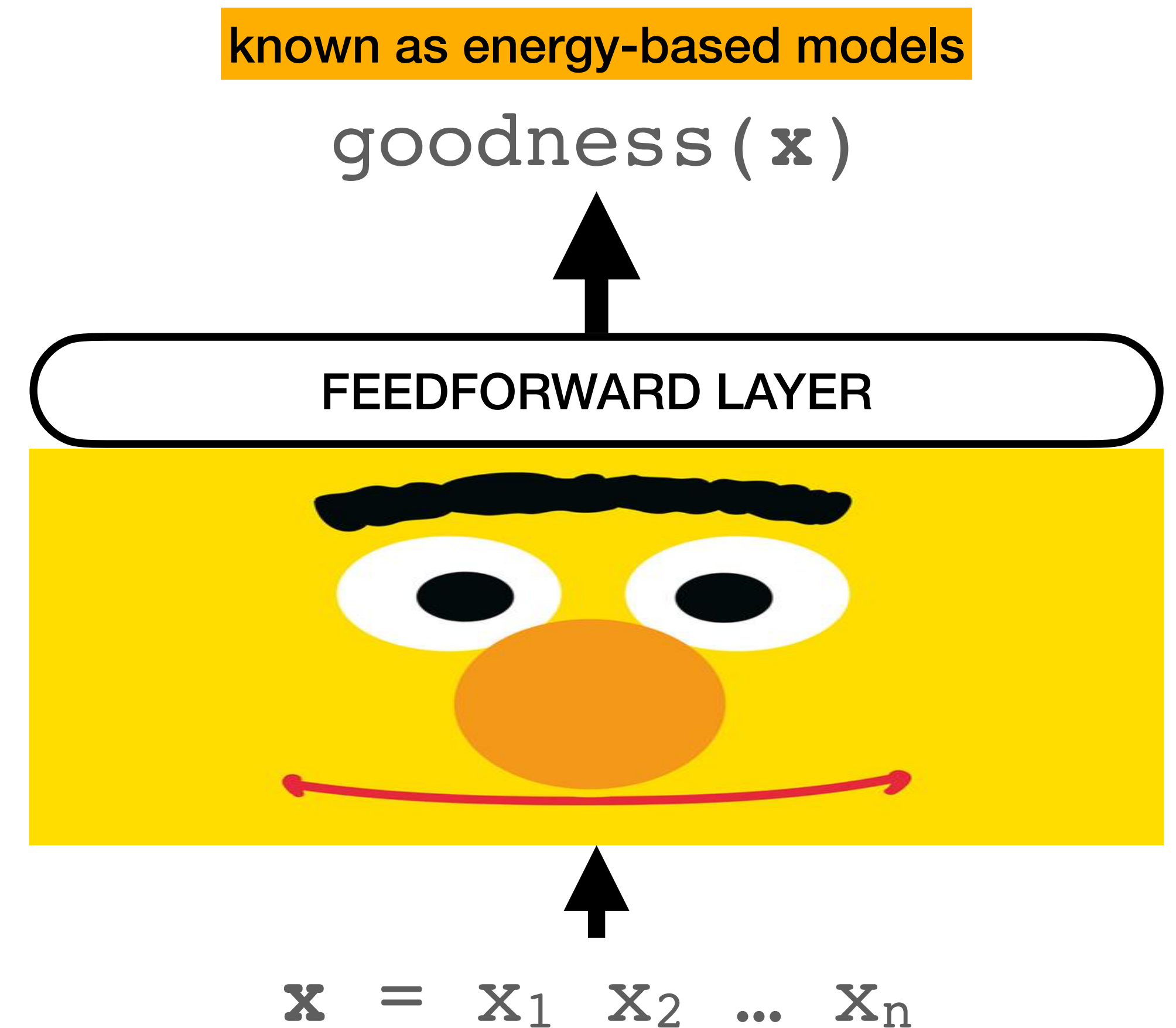
- **Good:**
Roses are red
- **Maybe:**
Roses are nosy
- **Bad:**
Roses queen sierra



Many NLP tasks are about scoring strings

we want to measure their goodness quantitatively with an NN

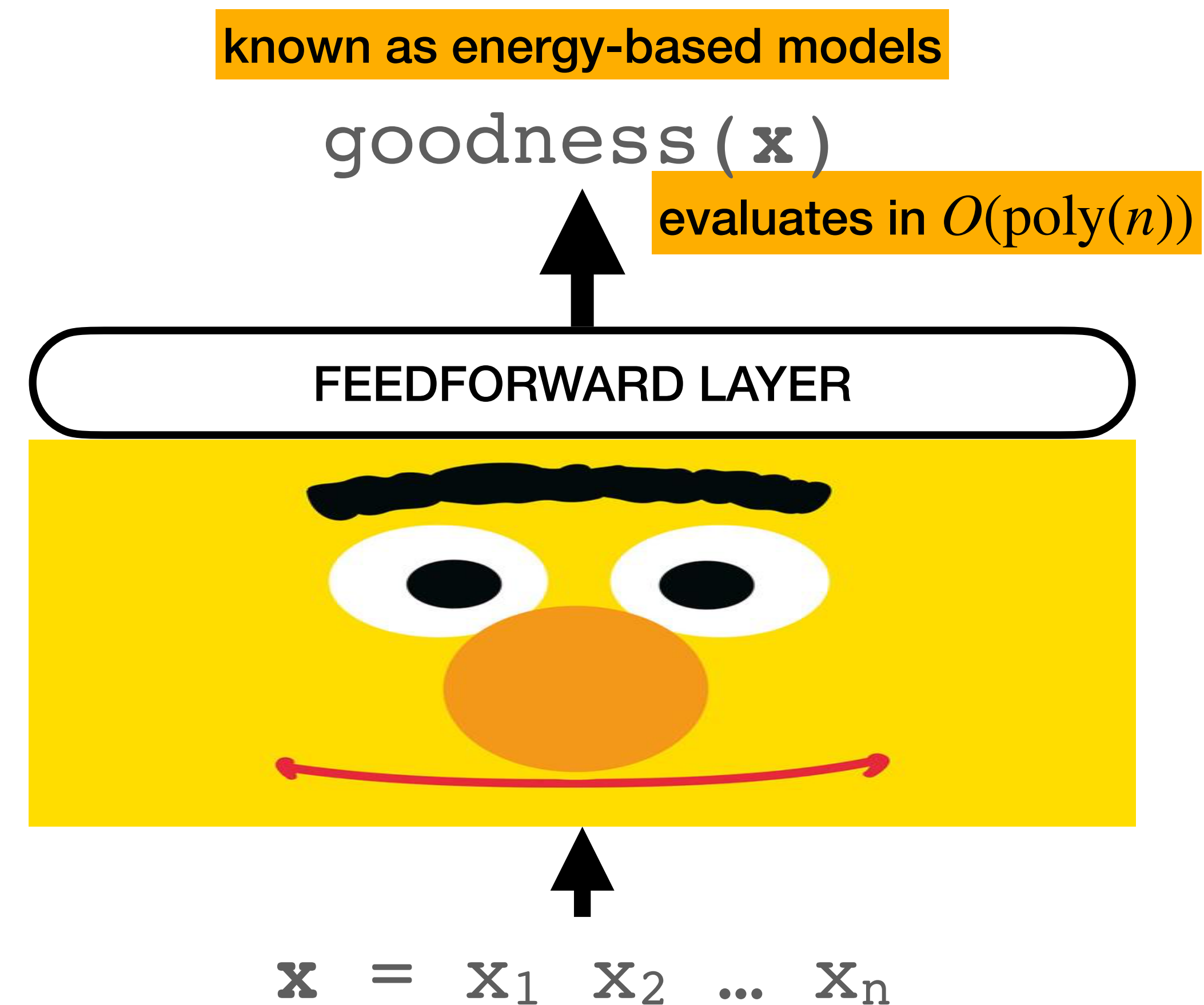
- $\text{goodness}(\text{"Roses are red"}) = 1000$
- $\text{goodness}(\text{"Roses are nosy"}) = 5$
- $\text{goodness}(\text{"Roses queen sierra"}) = 0.01$
- support of goodness:
 - set of strings whose goodness > 0



Many NLP tasks are about scoring strings

we want to measure their goodness quantitatively with an NN

- $\text{goodness}(\text{"Roses are red"}) = 1000$
- $\text{goodness}(\text{"Roses are nosy"}) = 5$
- $\text{goodness}(\text{"Roses queen sierra"}) = 0.01$
- support of goodness:
 - set of strings whose goodness > 0

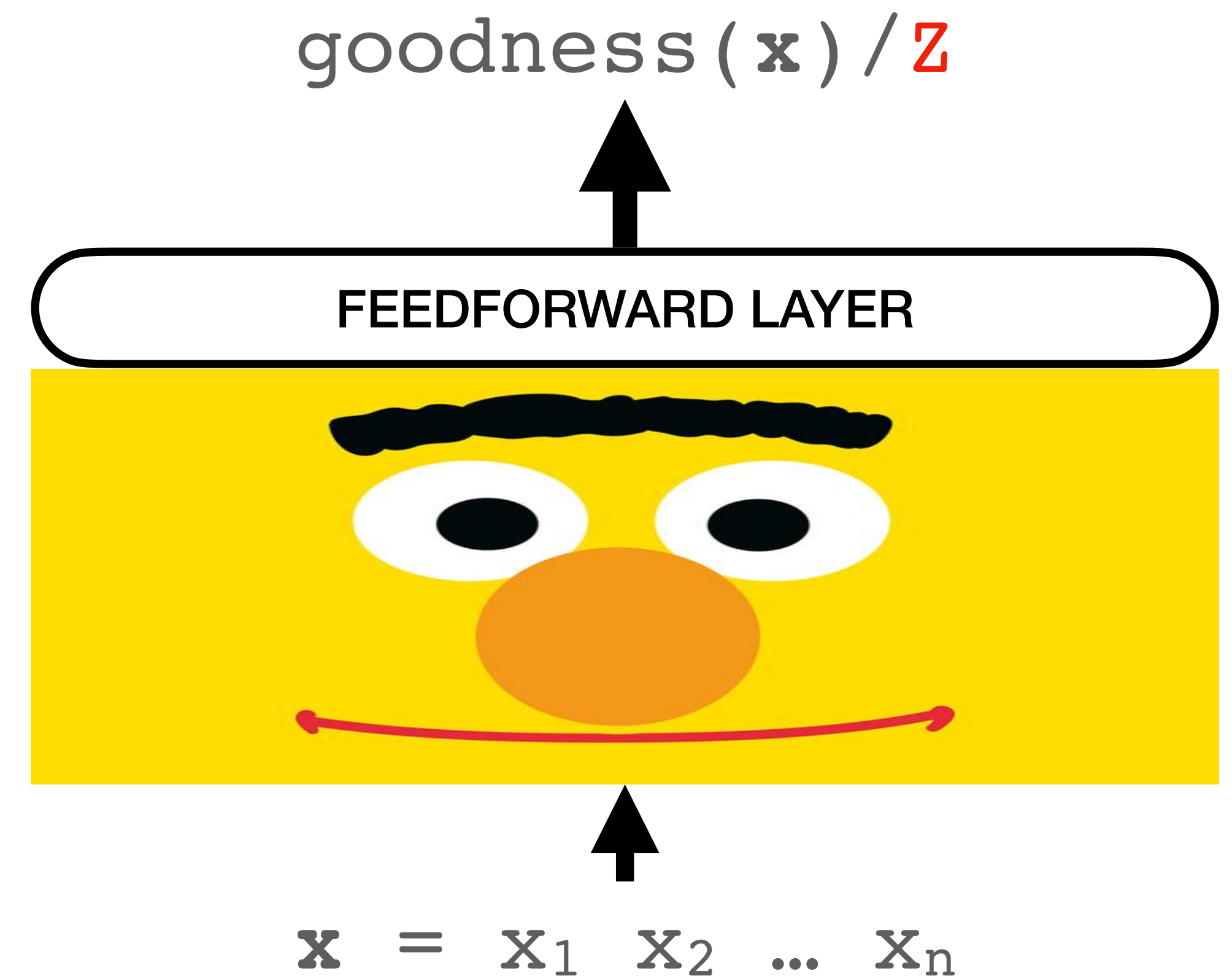


Many NLP tasks are about scoring strings

we want to measure their goodness on a scale between 0 and 1

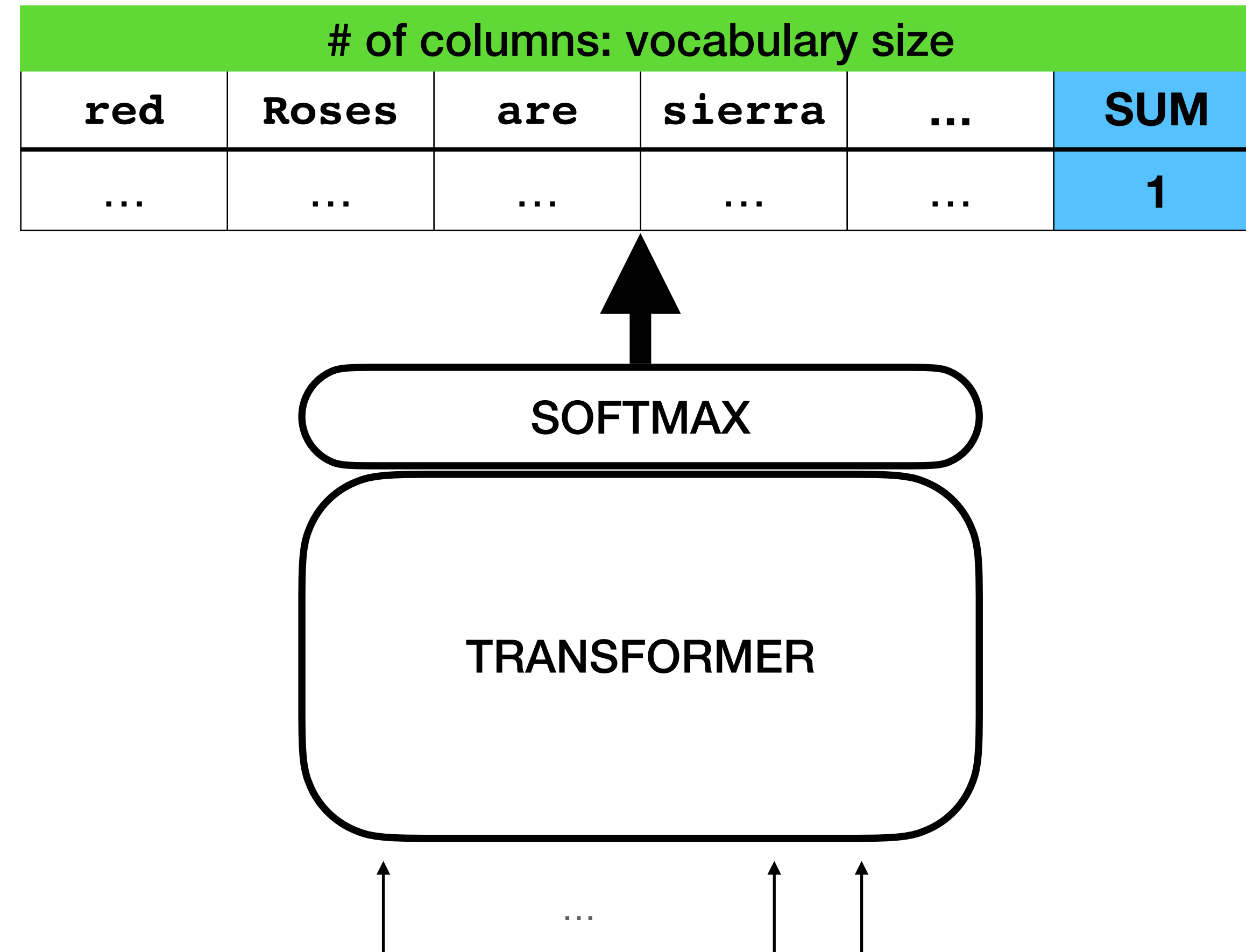
- $Z = \text{goodness}(\text{"Roses are red"})$
+ $\text{goodness}(\text{"Roses are nosy"})$
+ $\text{goodness}(\text{"Roses queen sierra"})$
+ ...
- intractable!

exponentially (or even infinitely) many columns!					
Roses are red	Roses are nosy	Roses queen	Las rosas son rojos	...	SUM
1000	5	0.01	0.00001	...	$1000+5+0.01+0.00001+\dots$



Autoregressive parametrization

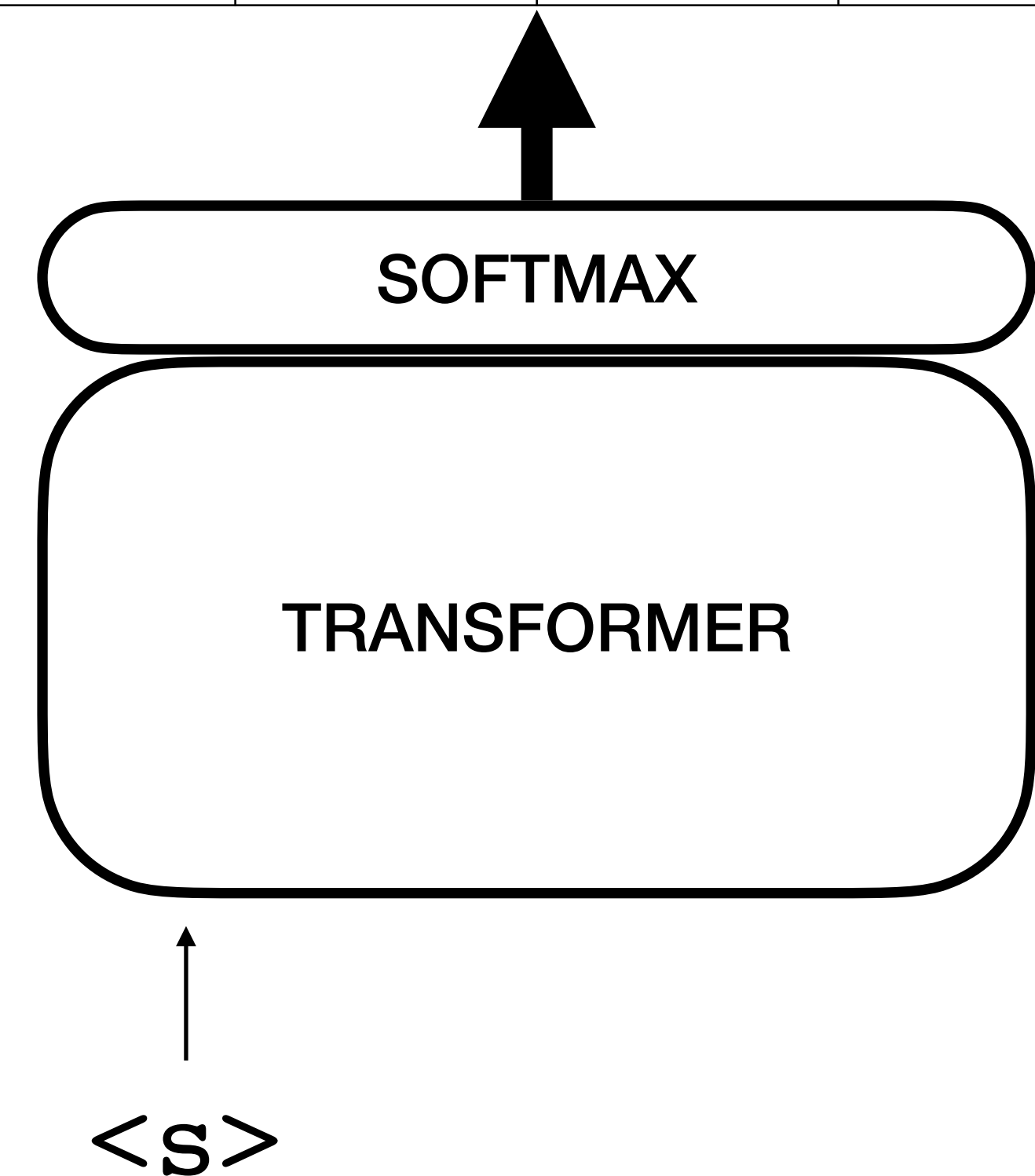
- `goodness("Roses are red") = ...`



Autoregressive parametrization

- goodness("Roses are red") =
 $p(\text{"Roses"}) \dots$
 $= 0.1 \dots$

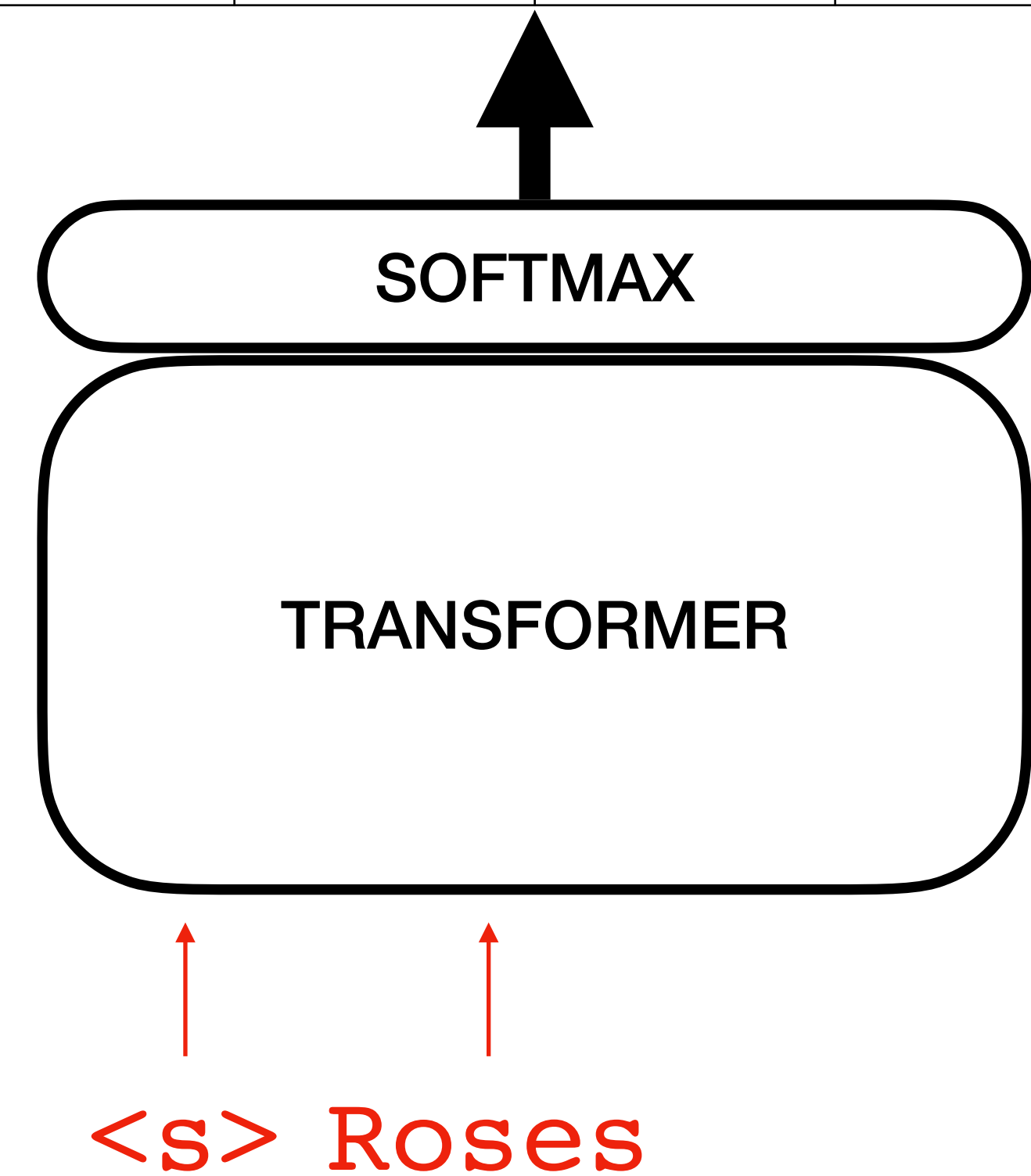
red	Roses	are	sierra	...	SUM
0.001	0.1	0.002	0.0001	...	1



Autoregressive parametrization

- goodness("Roses are red") =
 $p(\text{"Roses"}) * p(\text{"are"} | \text{"Roses"}) \dots$
 $= 0.1 * 0.3 \dots$

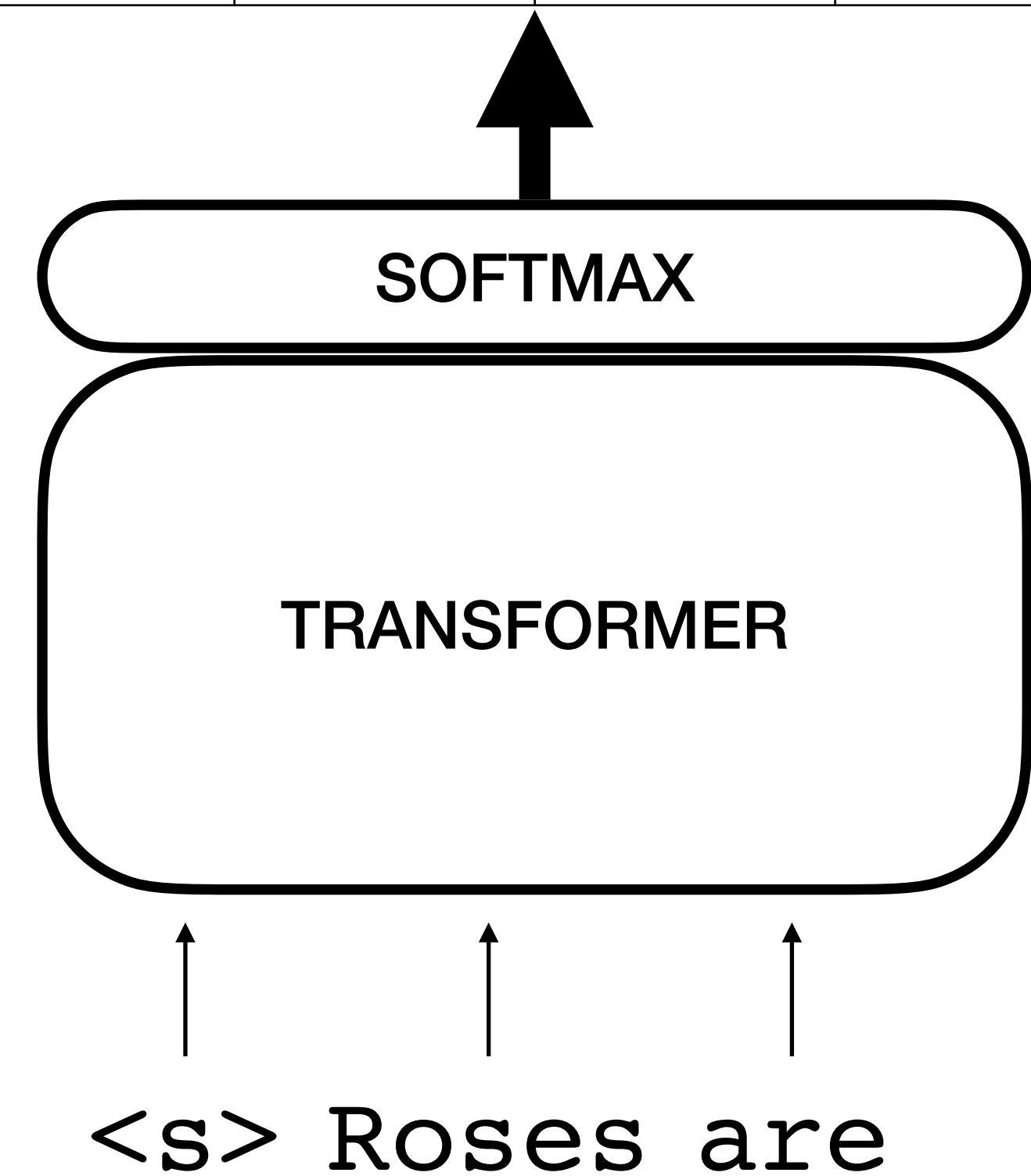
red	Roses	are	sierra	...	SUM
0.001	0.002	0.3	0.00002	...	1



Autoregressive parametrization

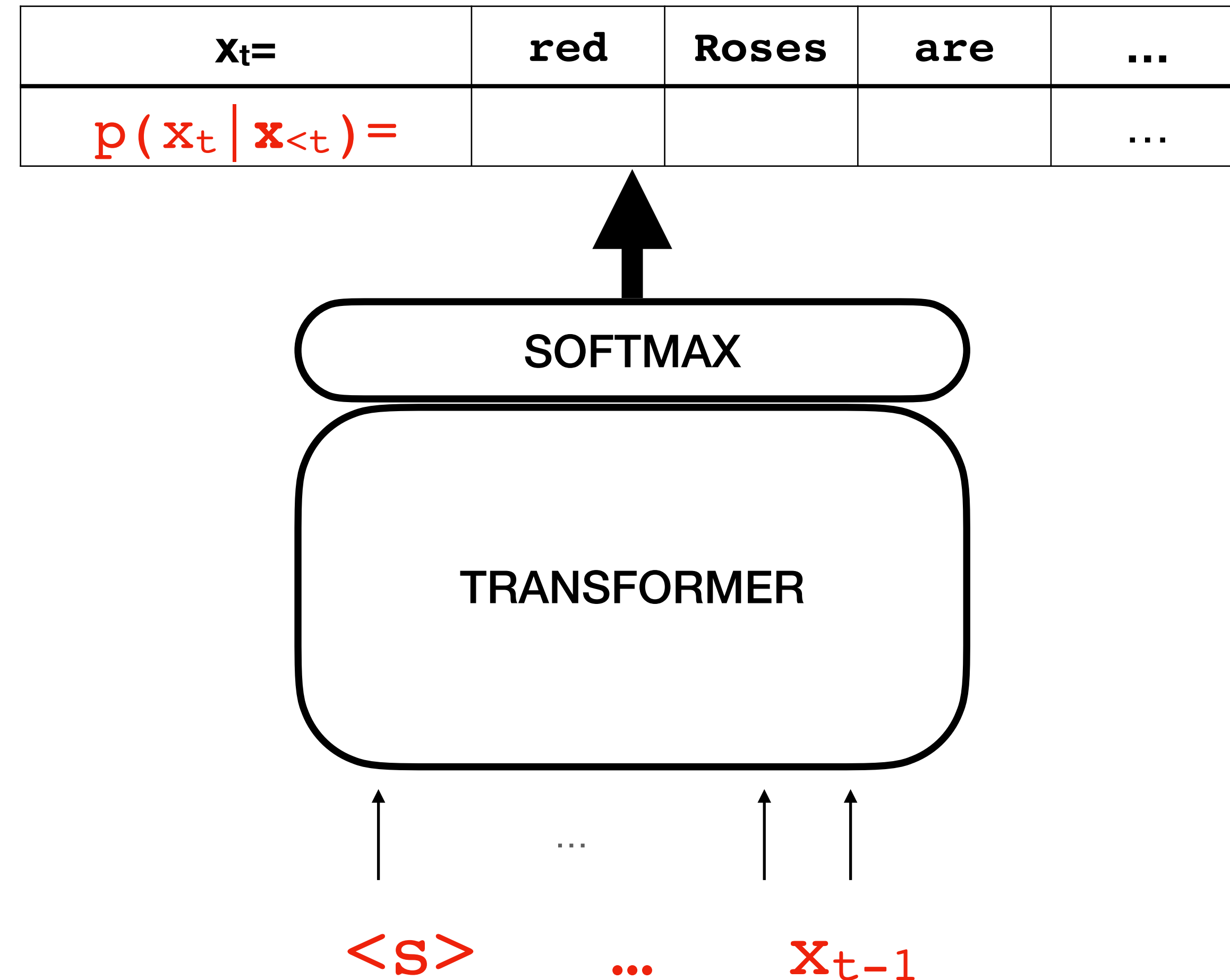
- goodness("Roses are red") =
 $p(\text{"Roses"}) * p(\text{"are"} | \text{"Roses"}) * p(\text{"red"} | \text{"Roses are"})$
 $= 0.1 * 0.3 * 0.05$

red	Roses	are	sierra	...	SUM
0.05	0.02	0.0007	0.0006	...	1



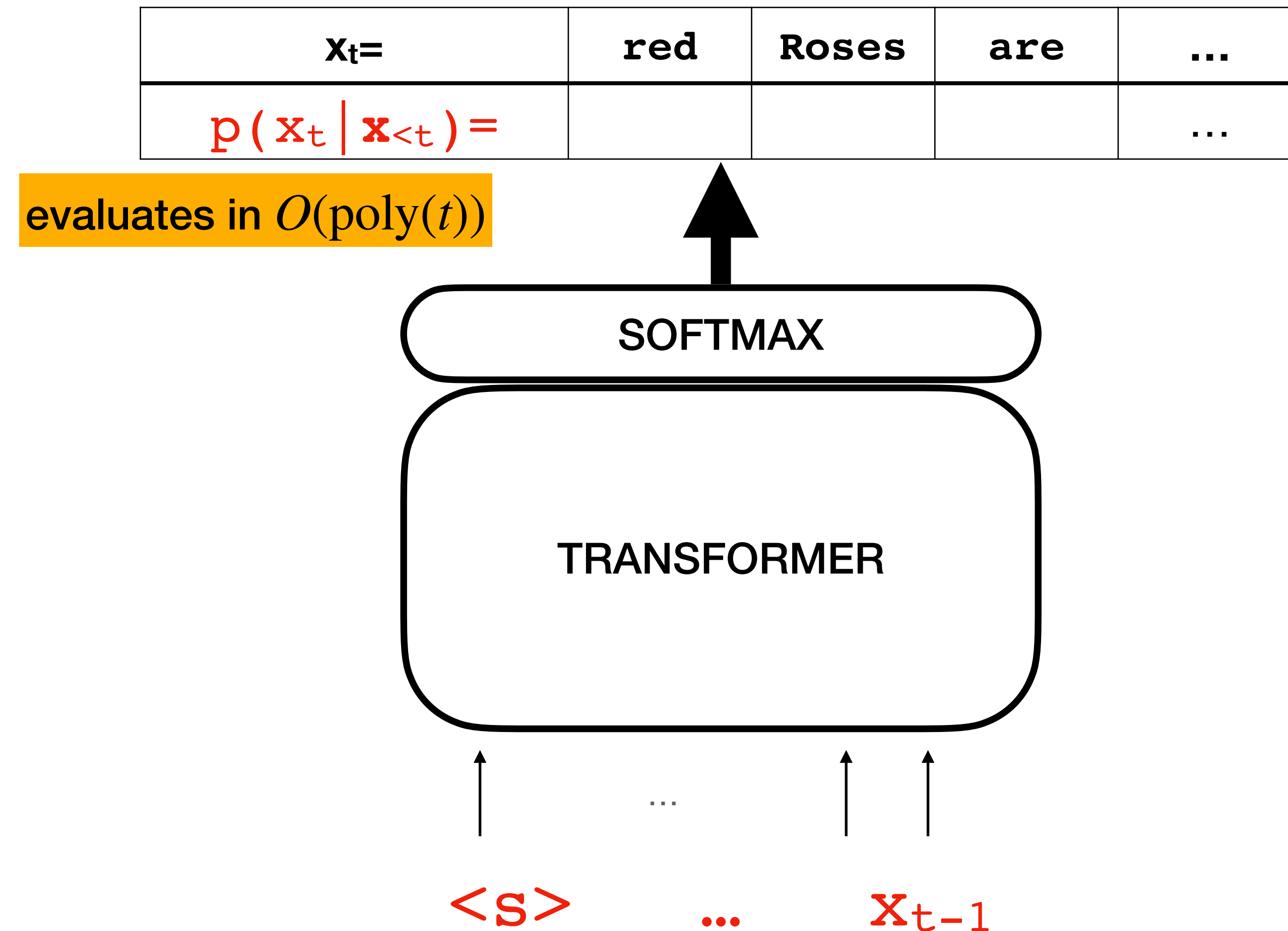
Autoregressive parametrization

- goodness("Roses are red")
=0.0015
- goodness("Roses are nosy")
=0.0000076
- goodness("Roses queen sierra")
=0.000000015



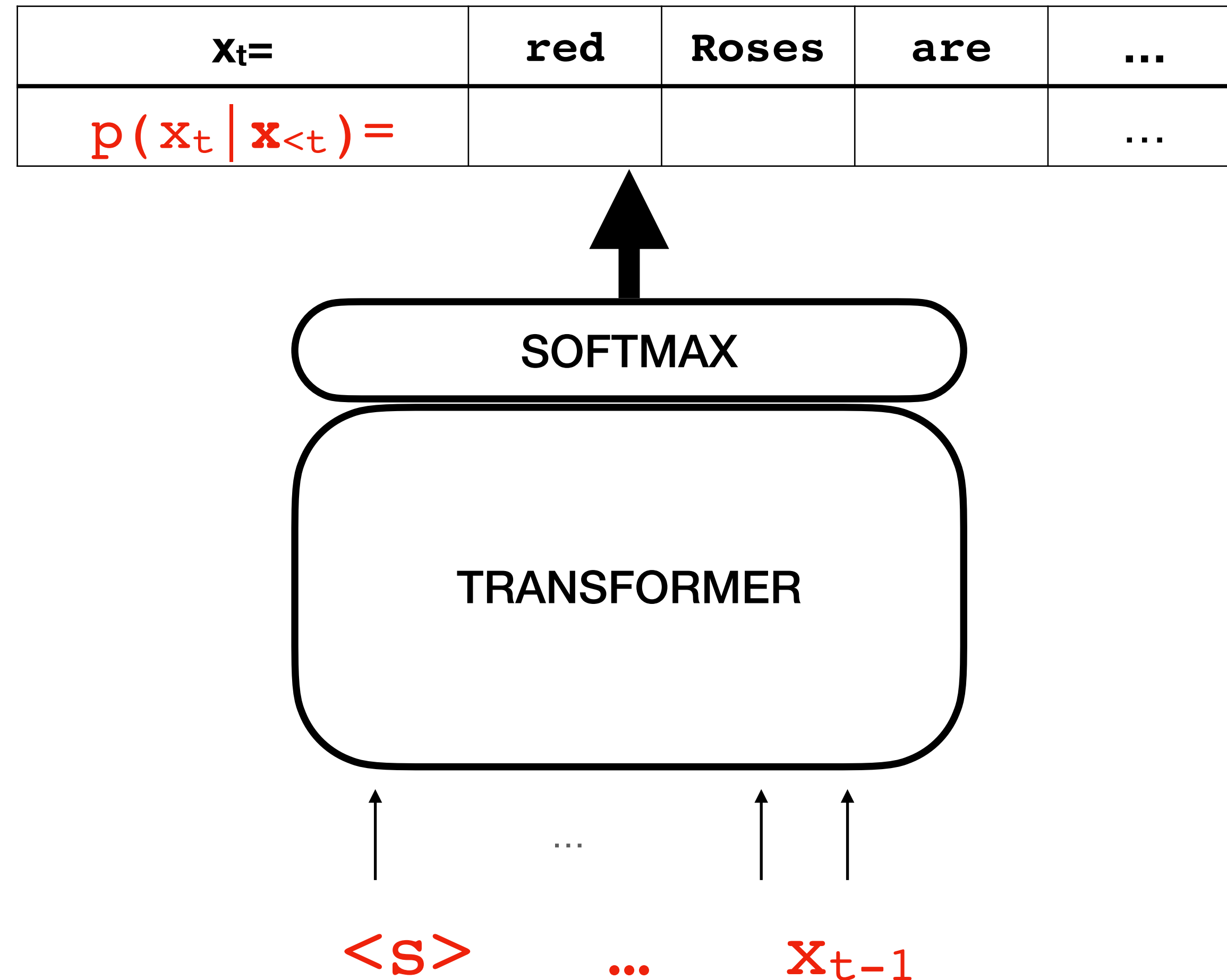
Autoregressive parametrization

- goodness("Roses are red")
=0.0015
- goodness("Roses are nosy")
=0.0000076
- goodness("Roses queen sierra")
=0.000000015

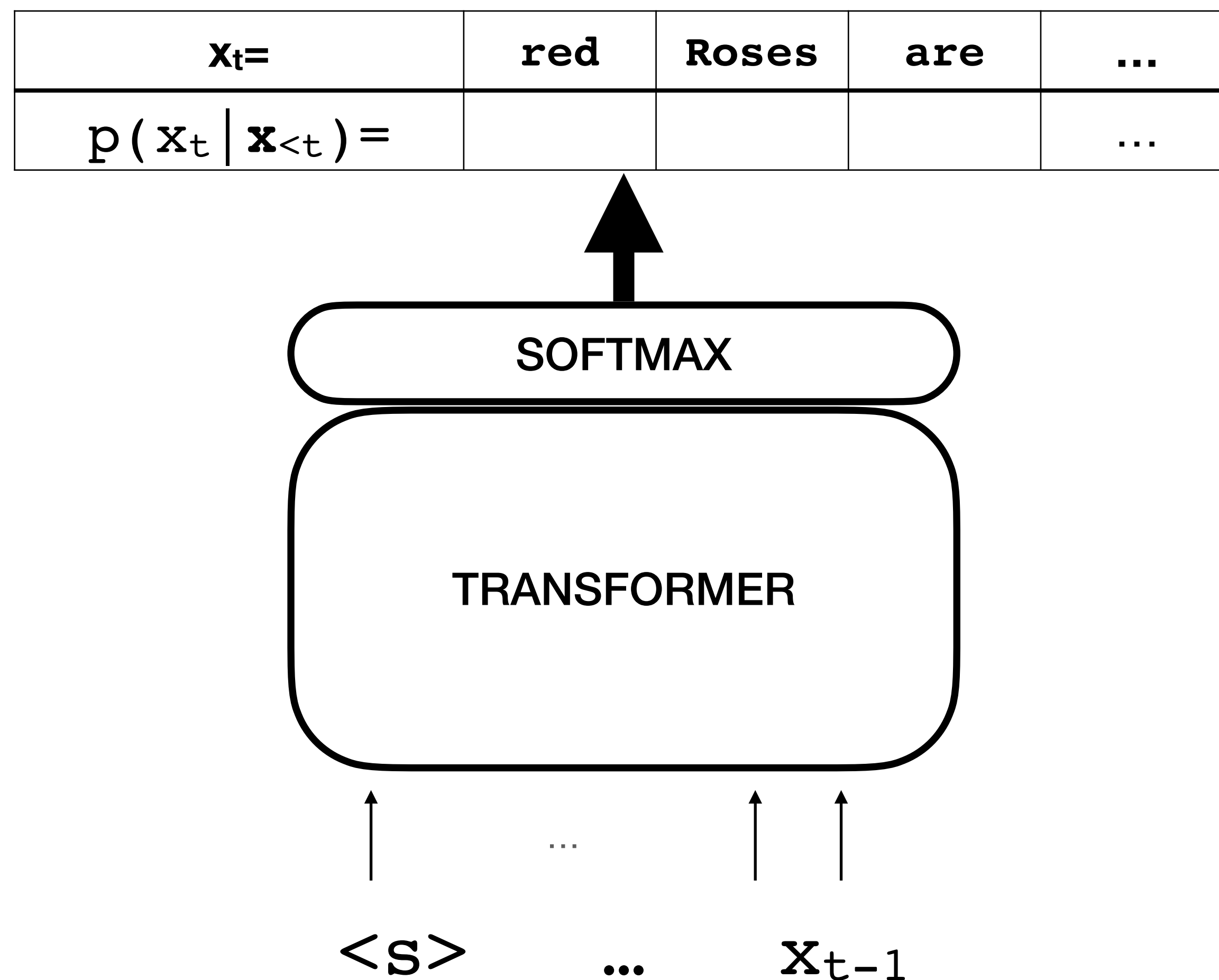
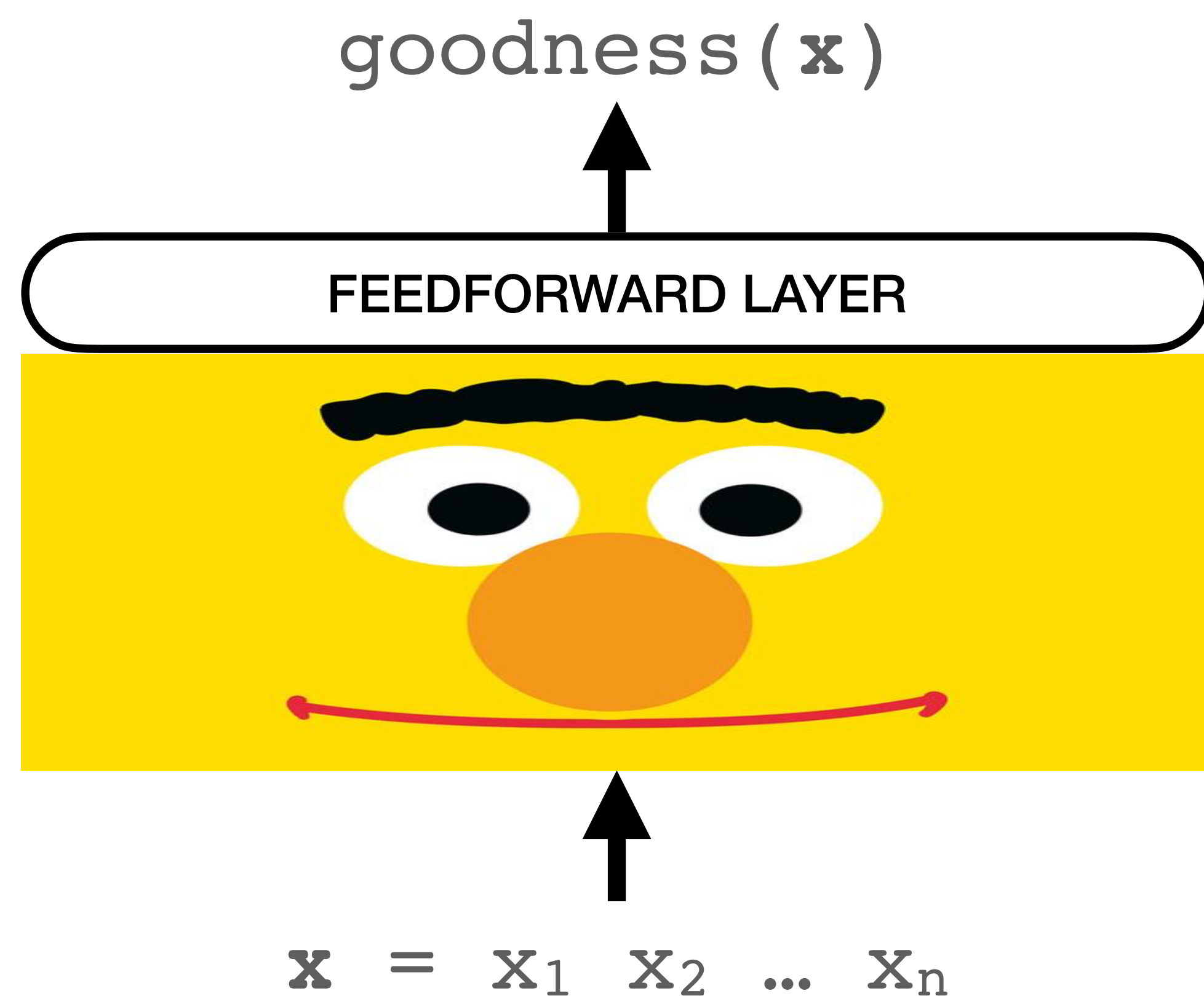


Autoregressive parametrization

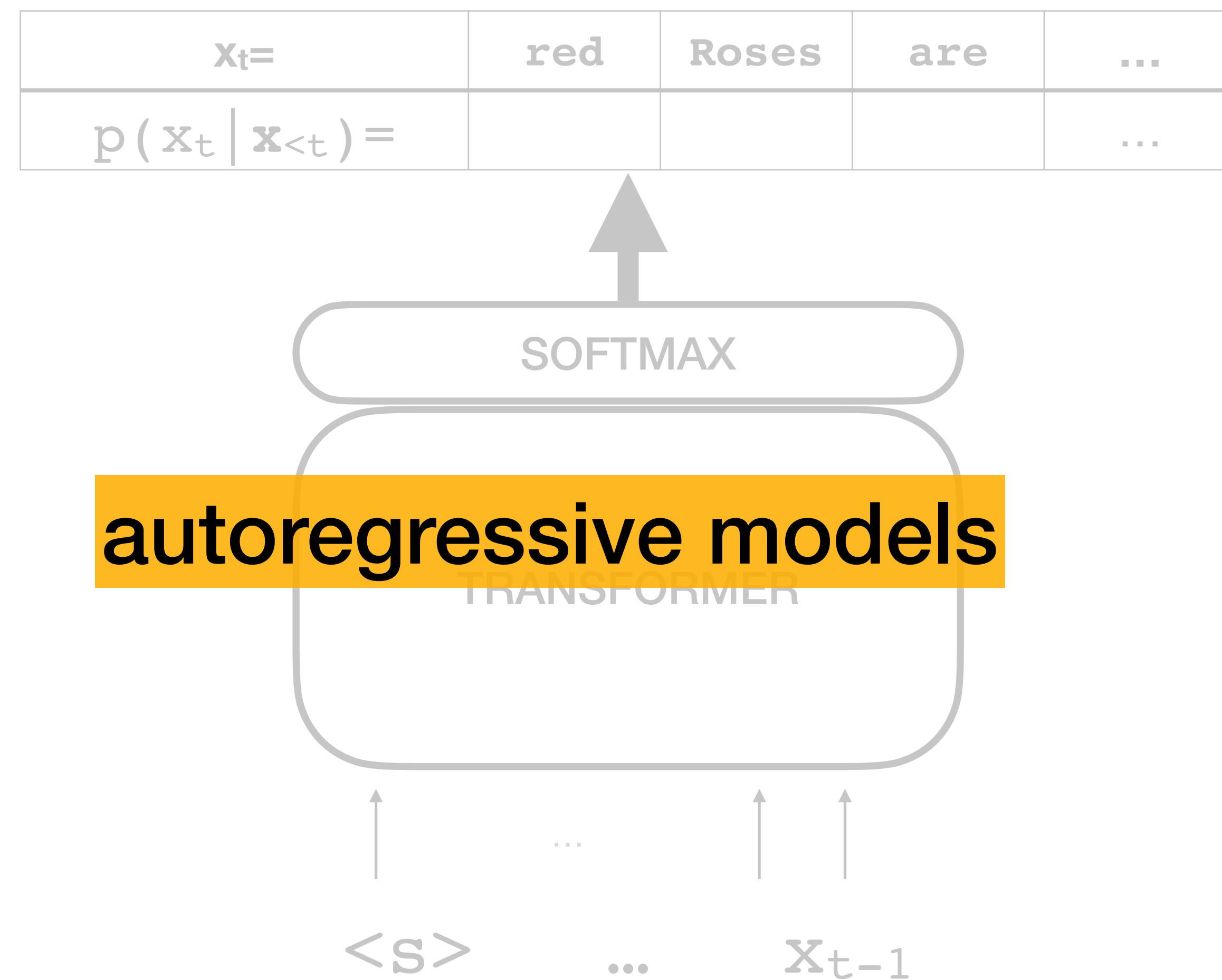
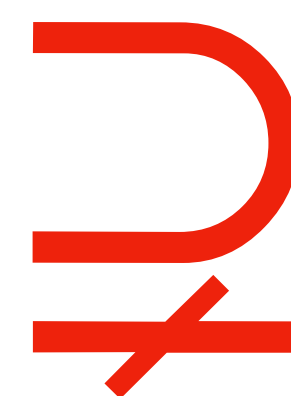
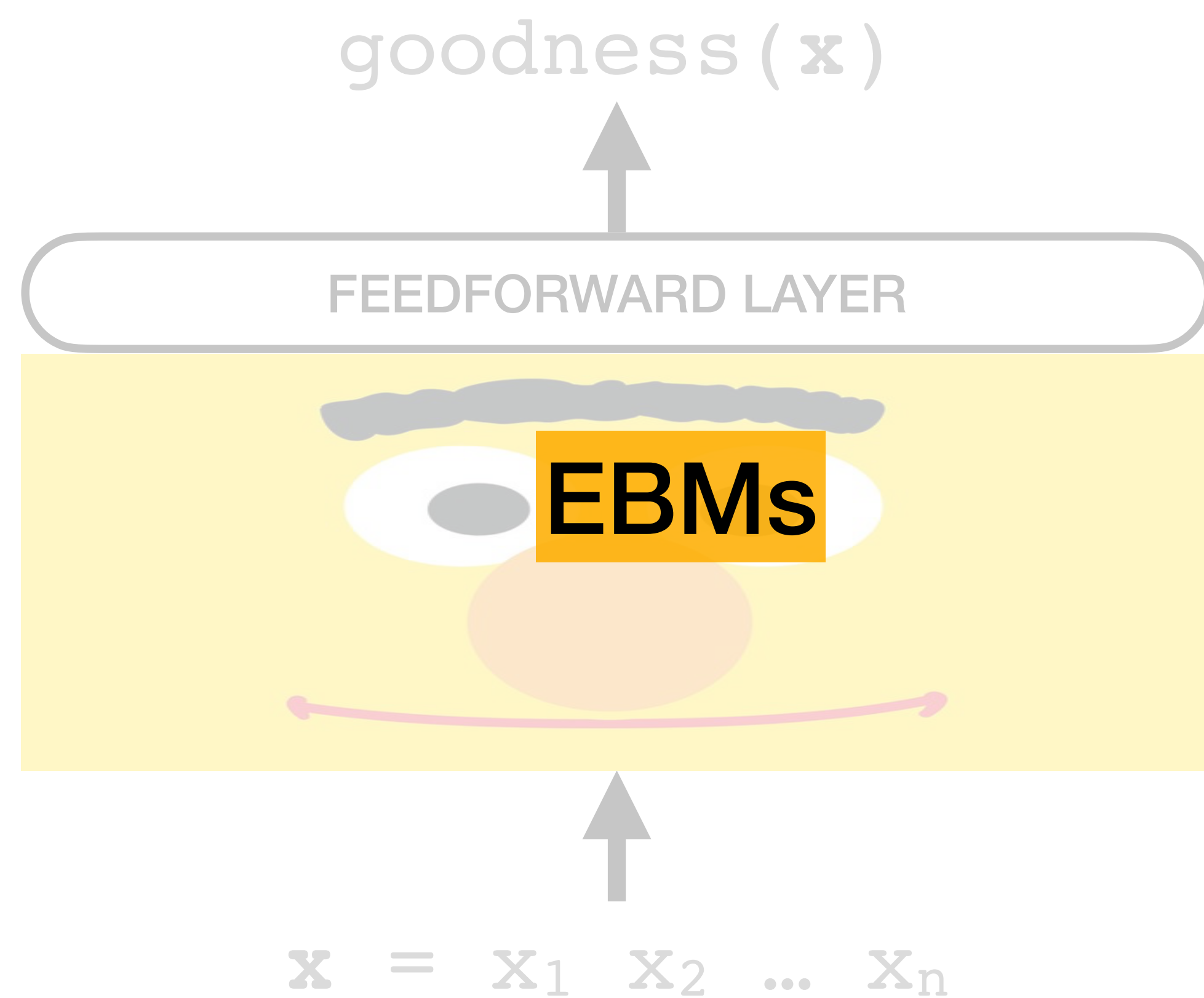
- Autoregressive models guarantee $Z = 1$
- $\text{goodness}(\text{"Roses are red"}) = 0.0015/Z = 0.0015$
- $\text{goodness}(\text{"Roses are nosy"}) = 0.00000076/Z = 0.00000076$
- $\text{goodness}(\text{"Roses queen sierra"}) = 0.0000000015/Z = 0.0000000015$



EBMs vs autoregressive models

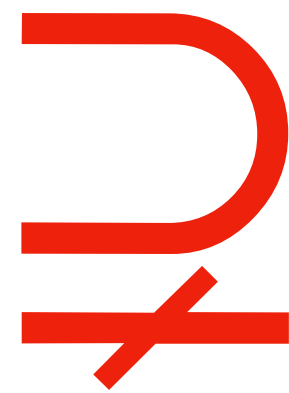


EBMs are more powerful!

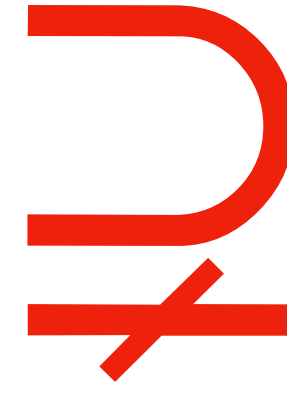


This work

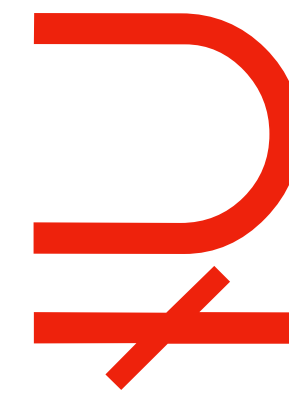
lookup
models



(autoregressive)
latent variable
models

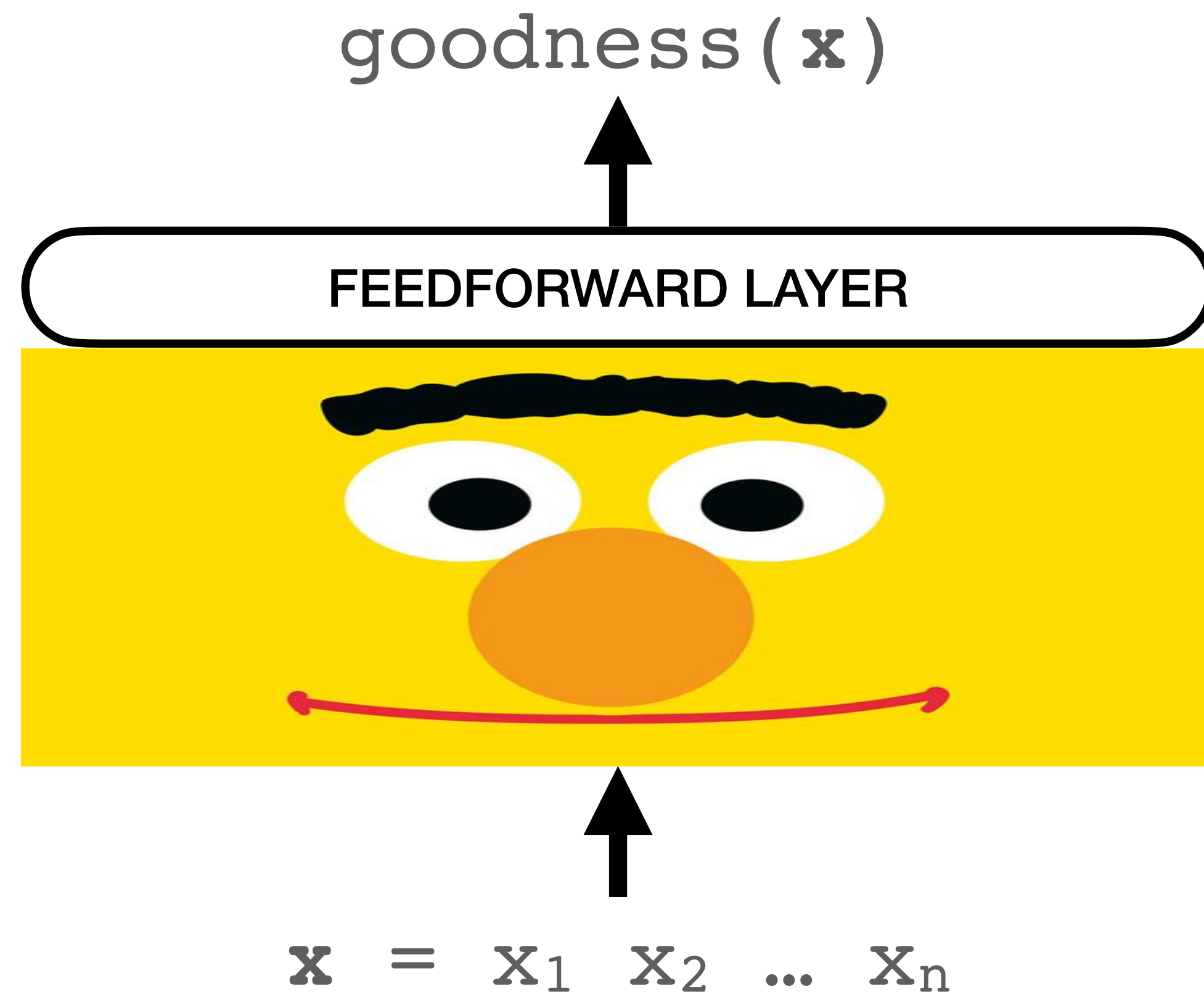


EBMs

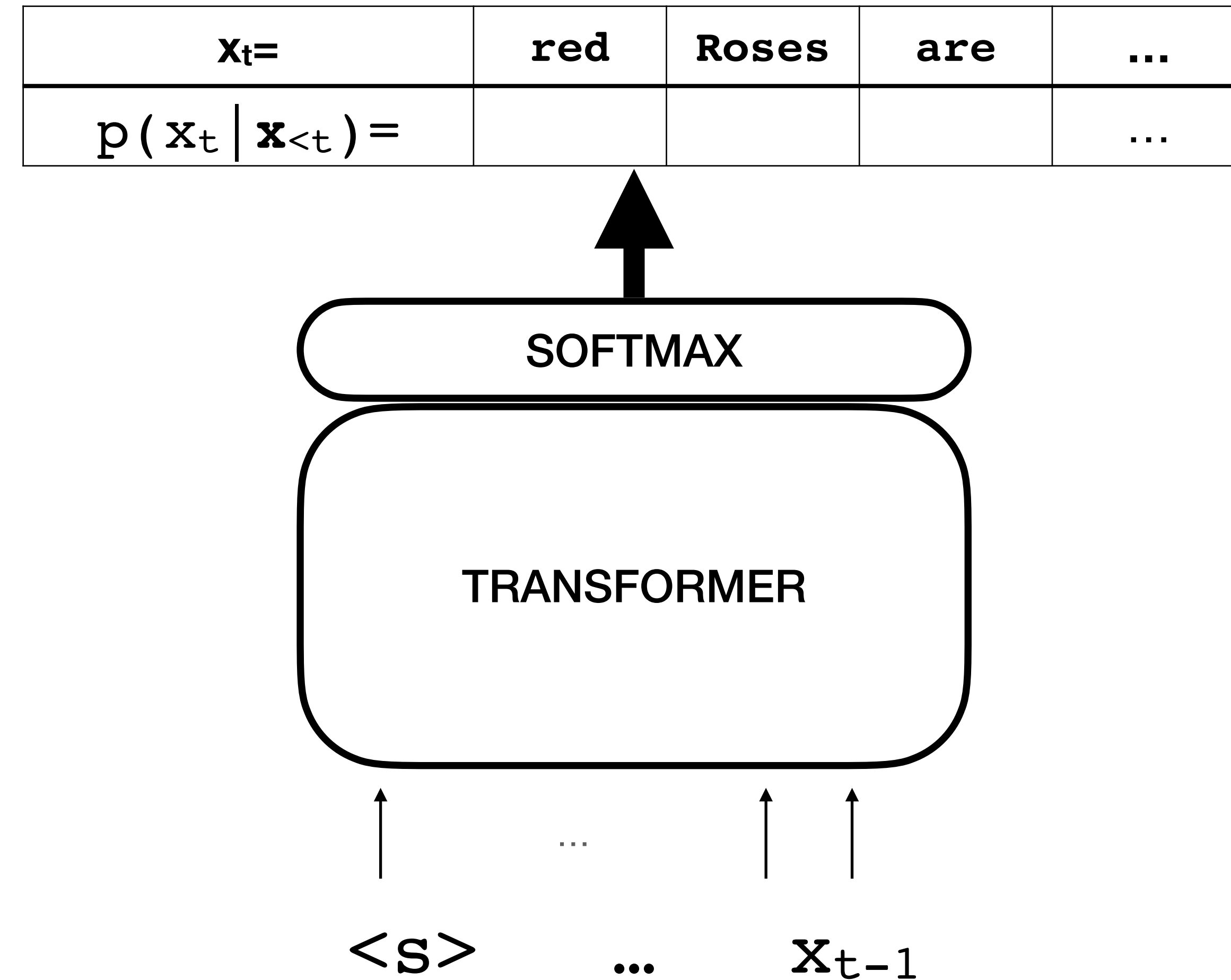


autoregressive
models

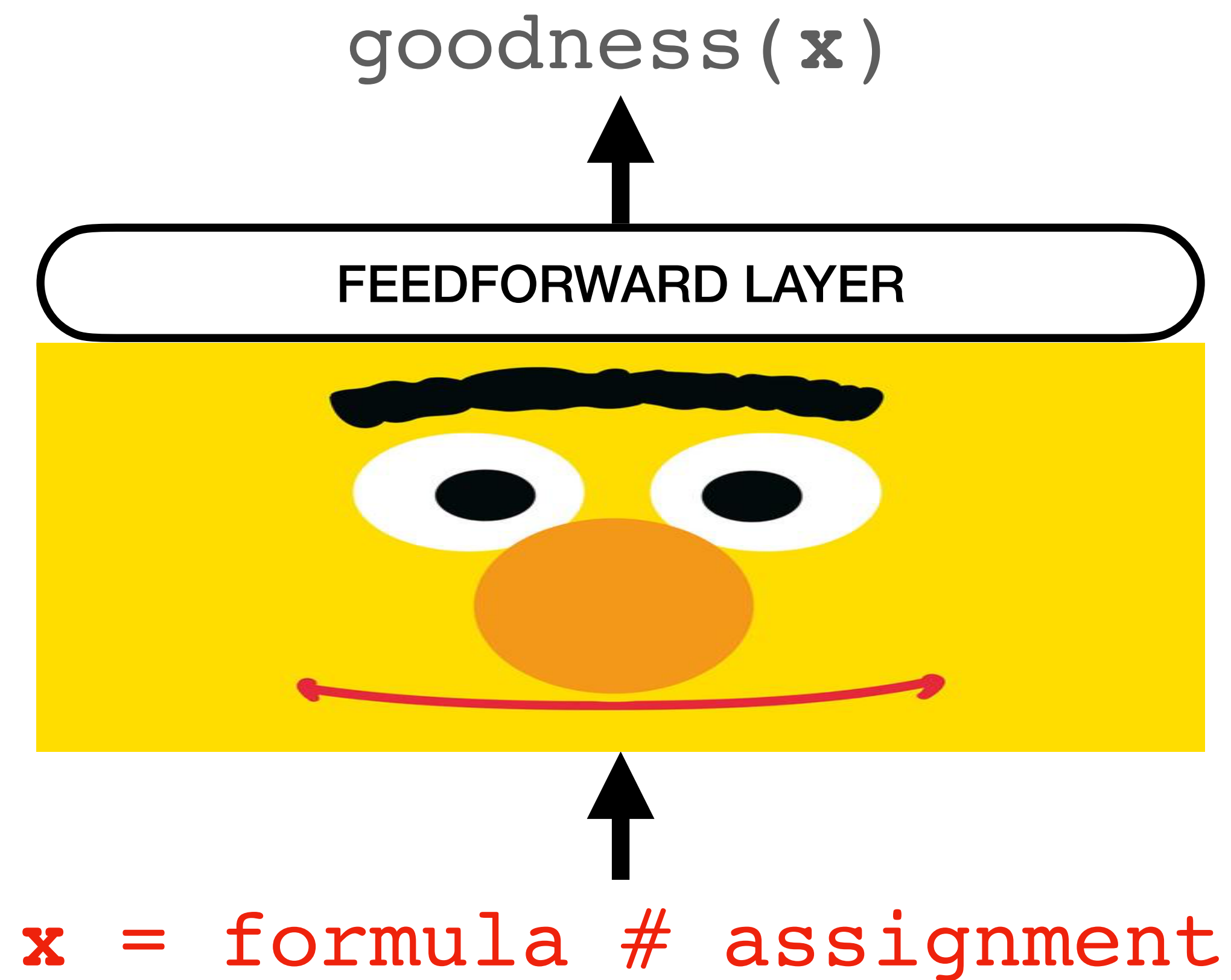
Why EBM $s \not\supseteq$ autoregressive models?



$\not\supseteq$

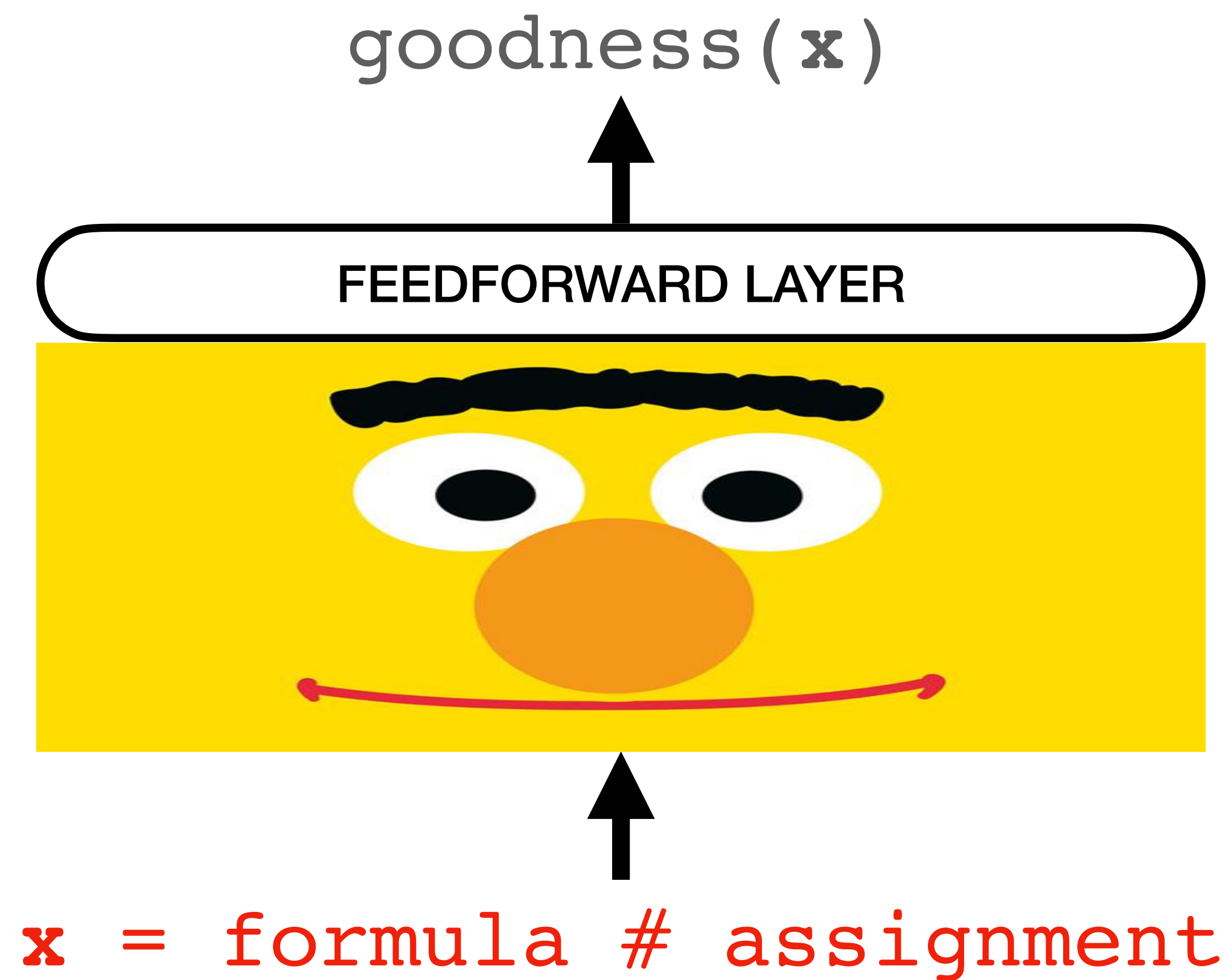


Why EBM $\not\supseteq$ autoregressive models?



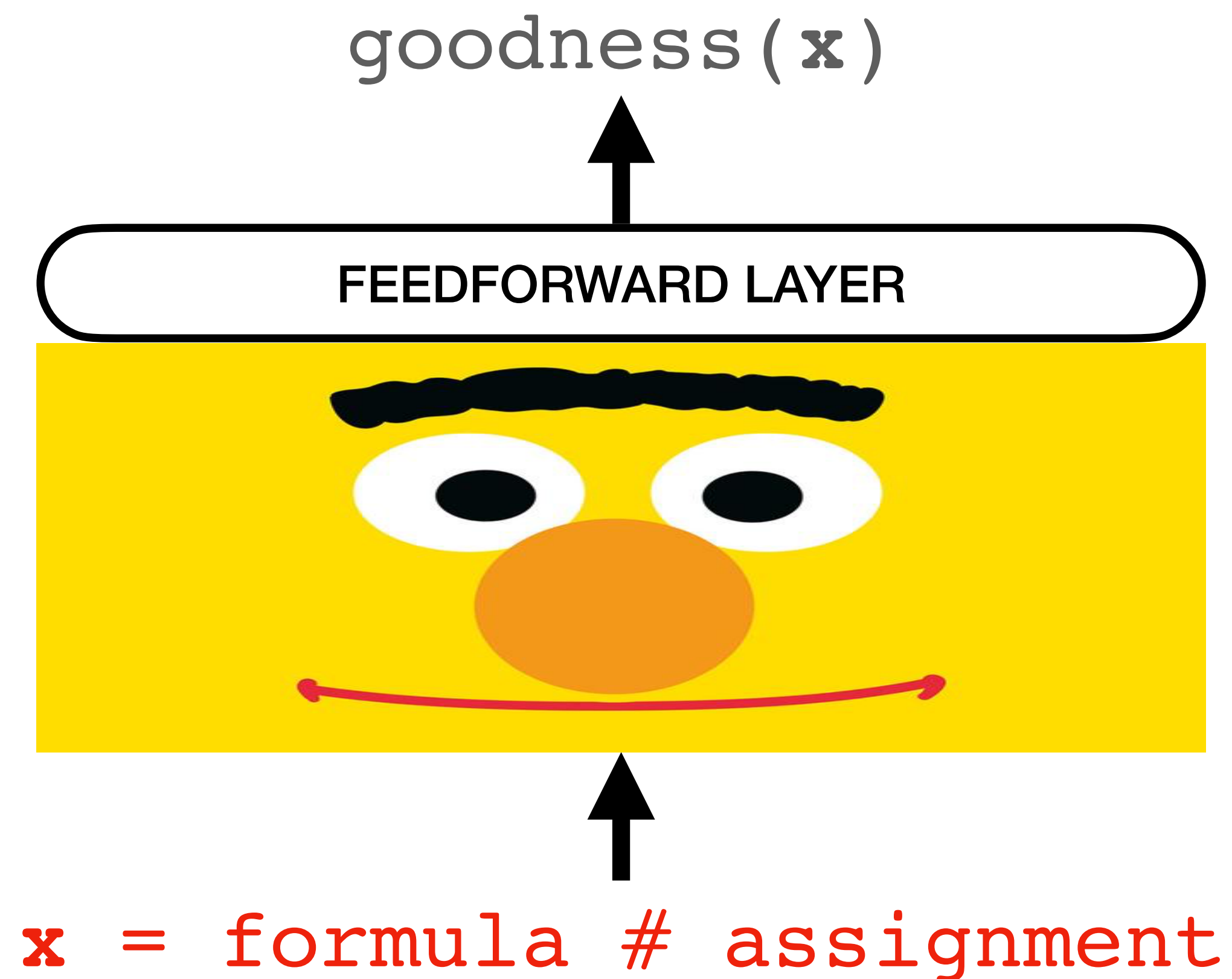
- formula:
 $(A_1 \text{ or not } A_2) \text{ and } (A_3)$
- assignment:
101

Why EBM $\not\supseteq$ autoregressive models?



- formula:
 $(A_1 \text{ or not } A_2) \text{ and } (A_3)$
- assignment:
101
- goodness(\mathbf{x})
 - > 0 if assignment satisfies formula
 - $=0$ otherwise

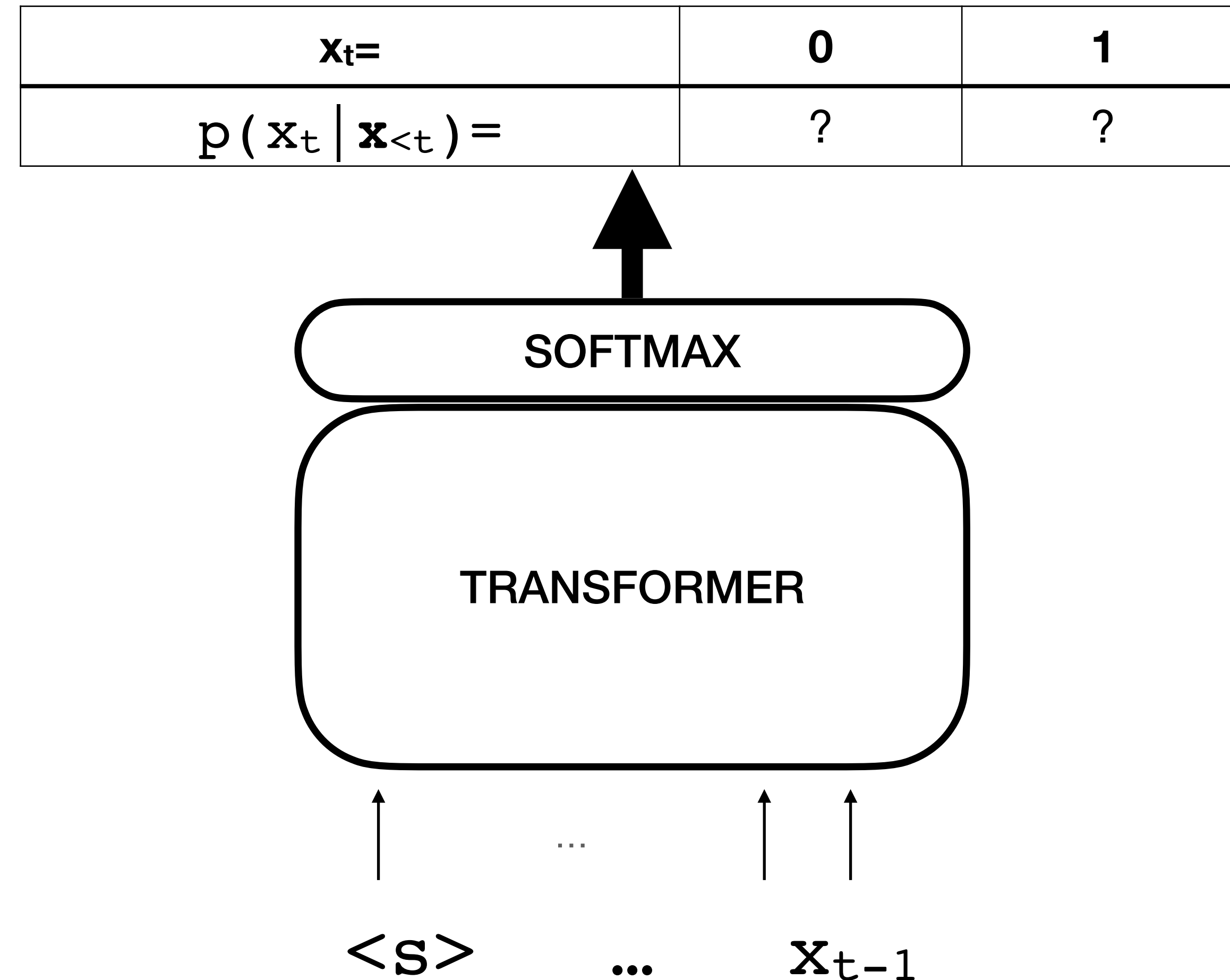
Why EBM $\not\supseteq$ autoregressive models?



- formula:
 $(A_1 \text{ or not } A_2) \text{ and } (A_3)$
- assignment:
1 0 1
- goodness(\mathbf{x})
 - > 0 if assignment satisfies formula
 - $=0$ otherwise
- goodness(\mathbf{x}) can be constructed as an RNN with size $O(|\mathbf{x}|^3)$

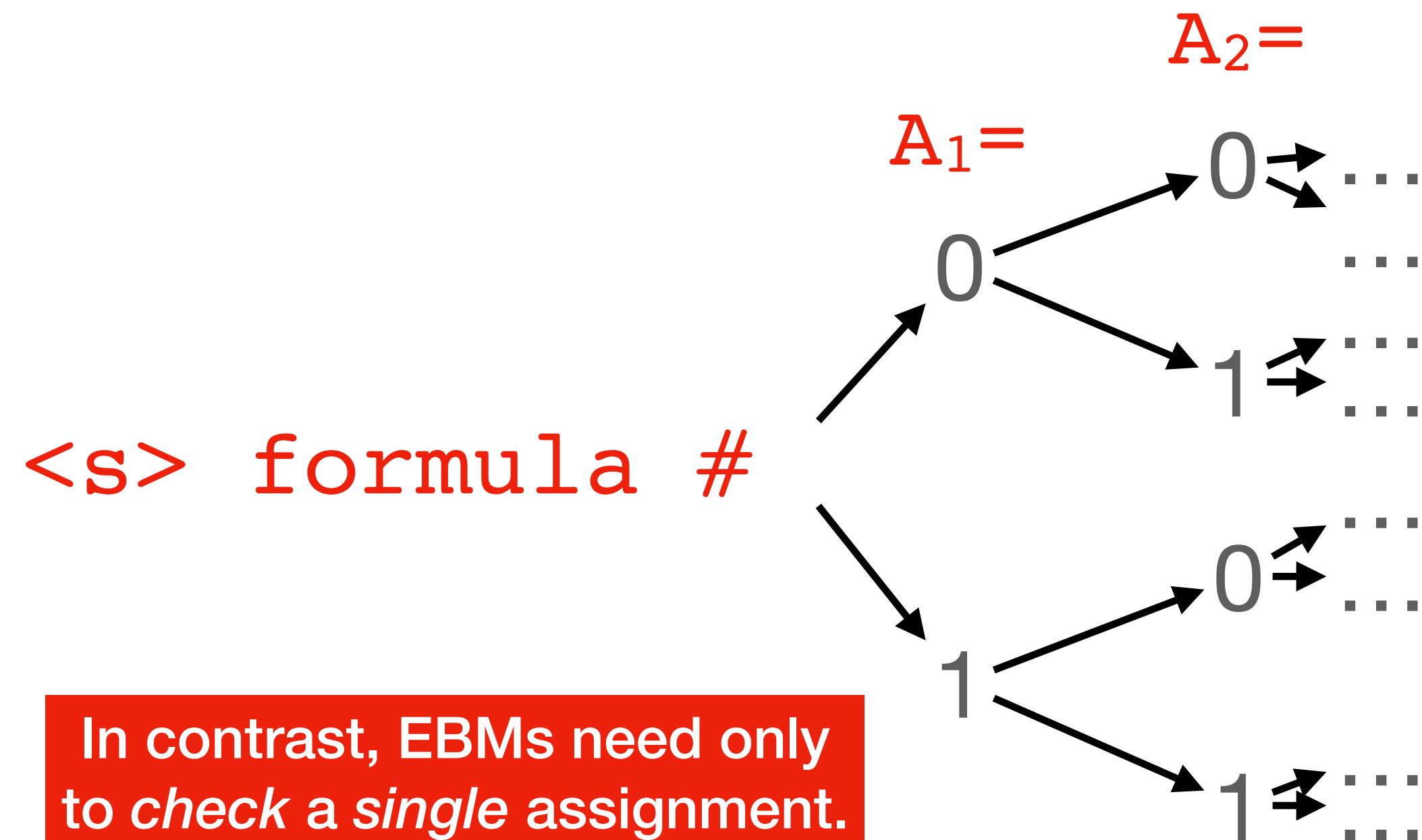
Why EBM_s $\not\supseteq$ autoregressive models?

- Now let's look at autoregressive models.
- Can we implement goodness(**x**) using an autoregressive model?

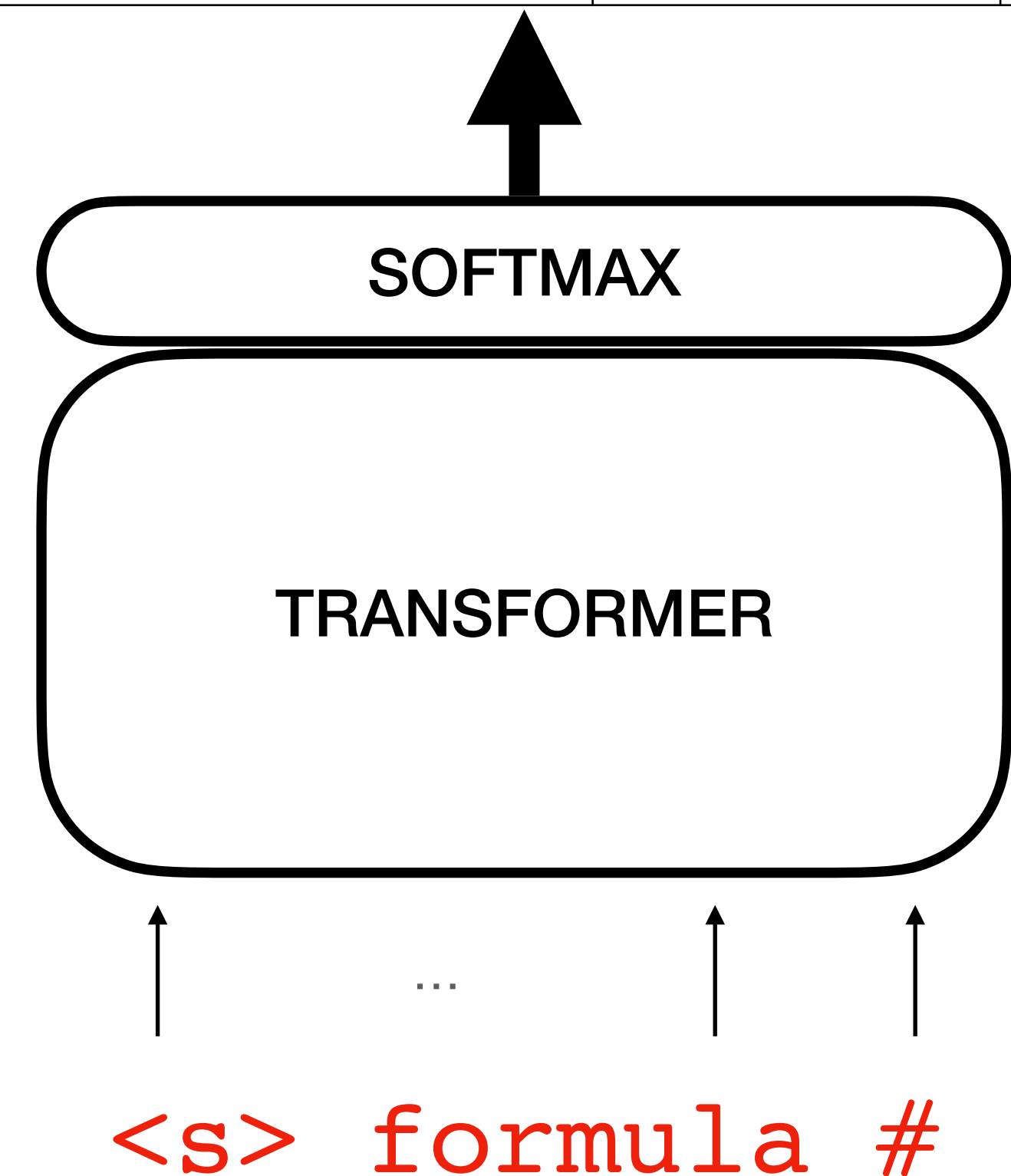


Why EBM $\not\supseteq$ autoregressive models?

- Computing the first token right after $\langle s \rangle$ formula # is as hard as determining if formula is satisfiable...

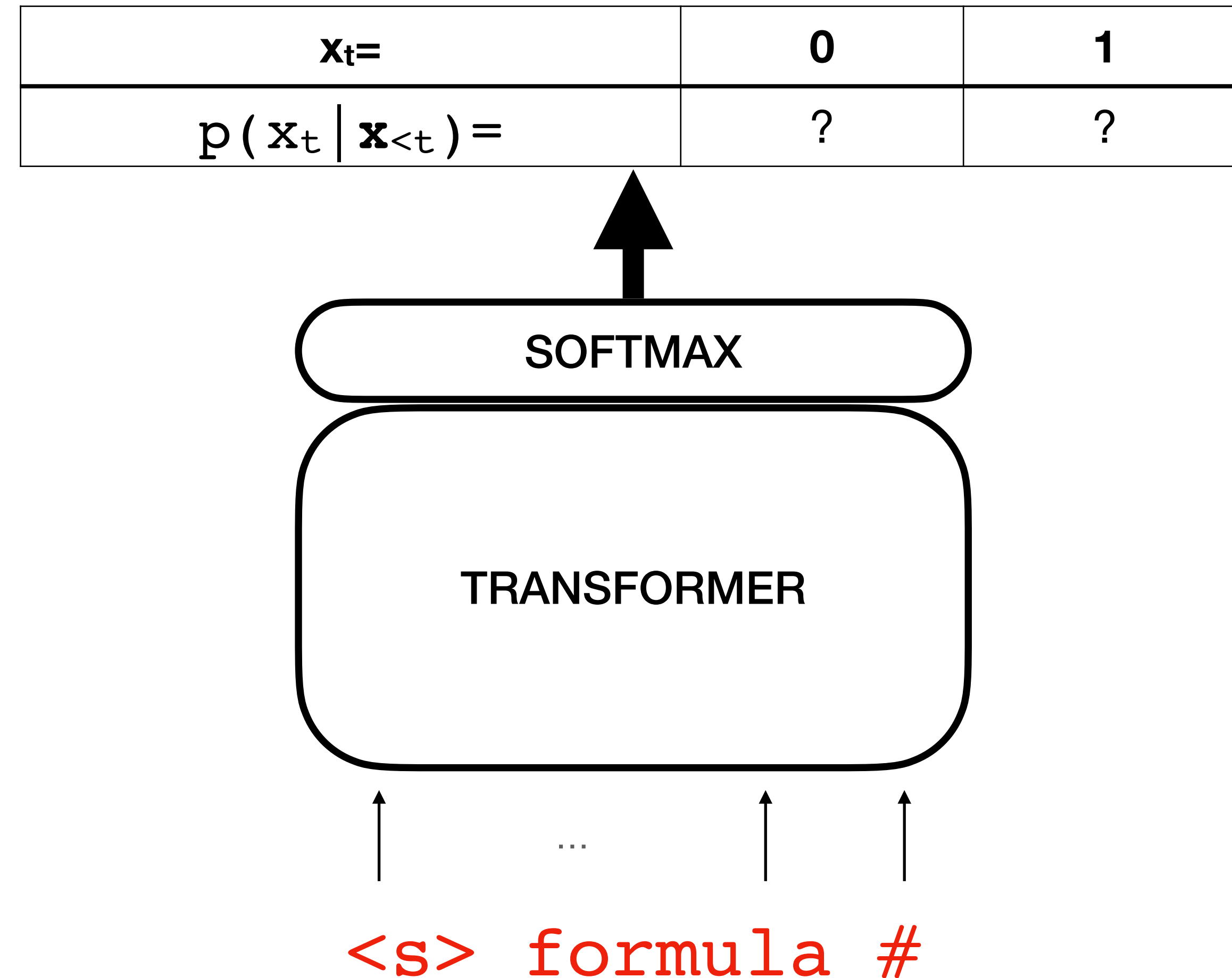


$x_t =$	0	1
$p(x_t \mathbf{x}_{<t}) =$?	?



Why EBM $\not\subseteq$ autoregressive models?

- Computing the first token right after $\langle s \rangle$ formula $\#$ is as hard as determining if formula is satisfiable...
 - which is NP-complete!
- Thus, **no** polynomial-time autoregressive model can model such distributions if $P \neq NP$.

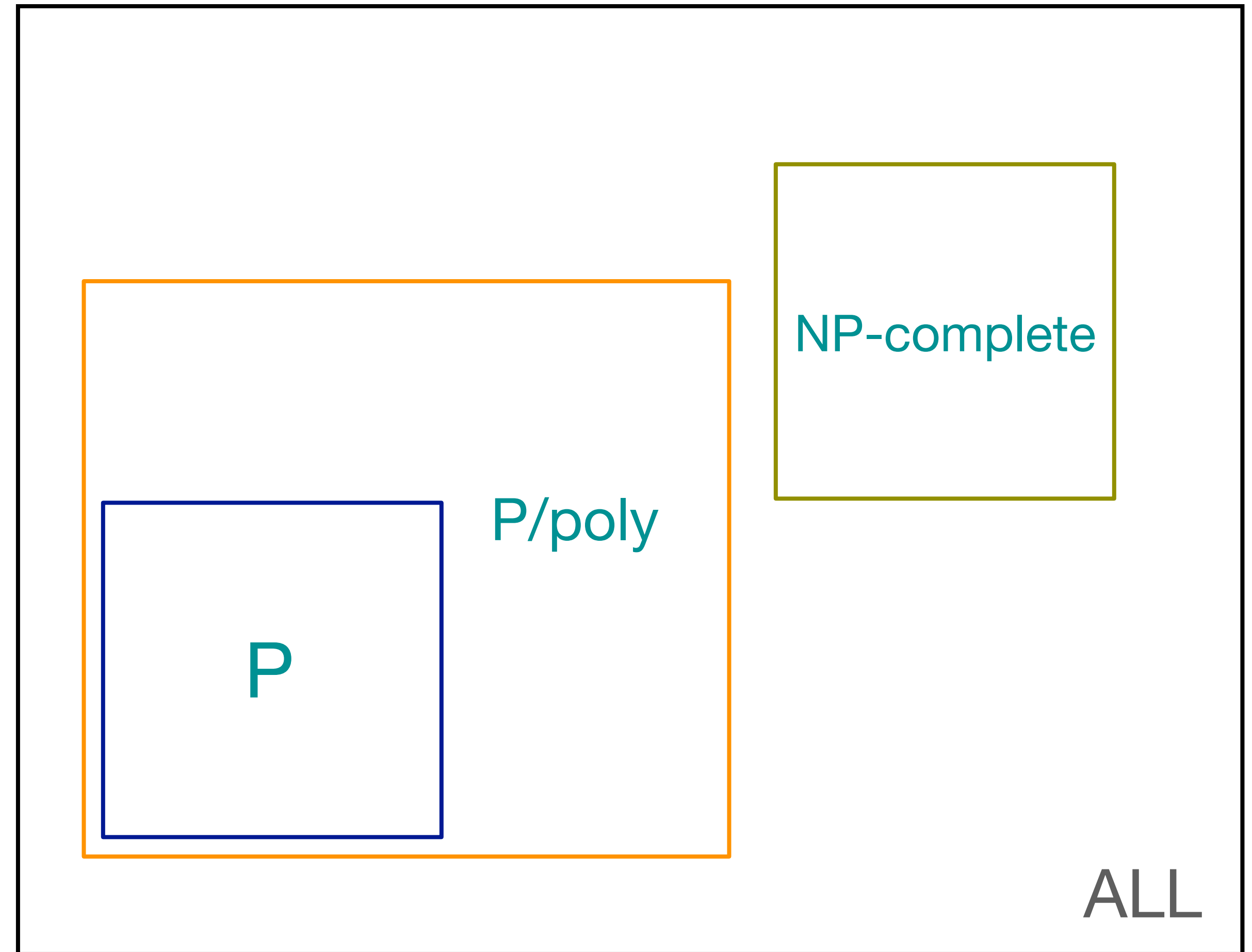


Actually it got worse

- There are distributions that can be captured by EBM.
But no autoregressive model can:
 - capture those distributions exactly (Theorem 1)
 - approximate well enough to get the same ranking of strings (Theorem 2)
 - approximate within any multiplicative factor (Theorem 4)
- Why should we care if we only need to model finite datasets?
 - in other words, we can always make the model larger to handle longer sequences (if smaller models don't work)...right?

Just make the model larger?

- If the model sizes *only* grow polynomially in sequence length, they belong in the P/poly class.
- It is widely believed $\text{NP} \not\subseteq \text{P/poly}$.
- So the models must **grow superpolynomially larger** and/or **run superpolynomially longer** in sequence length, to model longer problems (since they are NP-hard).
 - otherwise, the model simply won't fit even with access to an oracle in training!



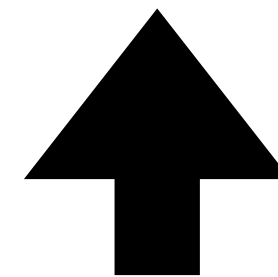
What if I am not interested in Boolean SAT problems?

- In general, autoregressive models cannot capture distributions over strings of the form `problem#solution`, where a `problem` is computationally hard to `solve`.
 - EBM's can capture such distributions
- Some CL/NLP `problems` are indeed computationally hard:
 - Parsing of many syntactic/semantic formalisms (e.g. AMR)
 - Propositional logic (NLI)
 - Optimality Theory
- Important linguistic regularities cannot be captured by autoregressive models!
- We use propositional logic generating a Star Wars movie script as an example.





A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



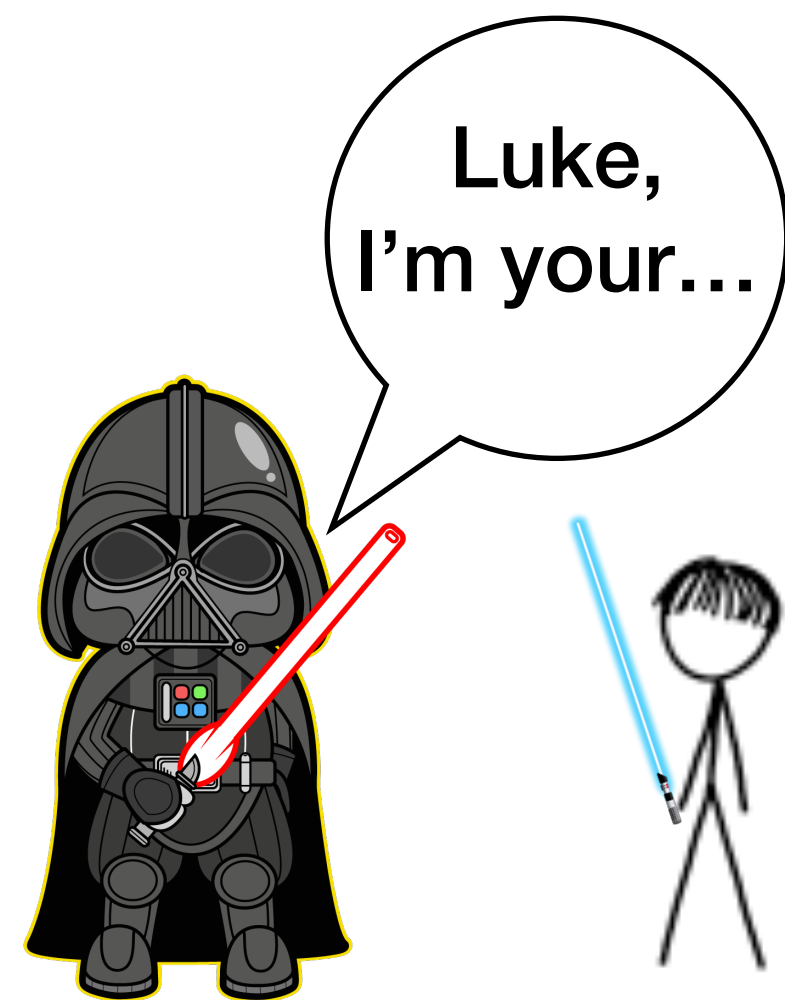
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

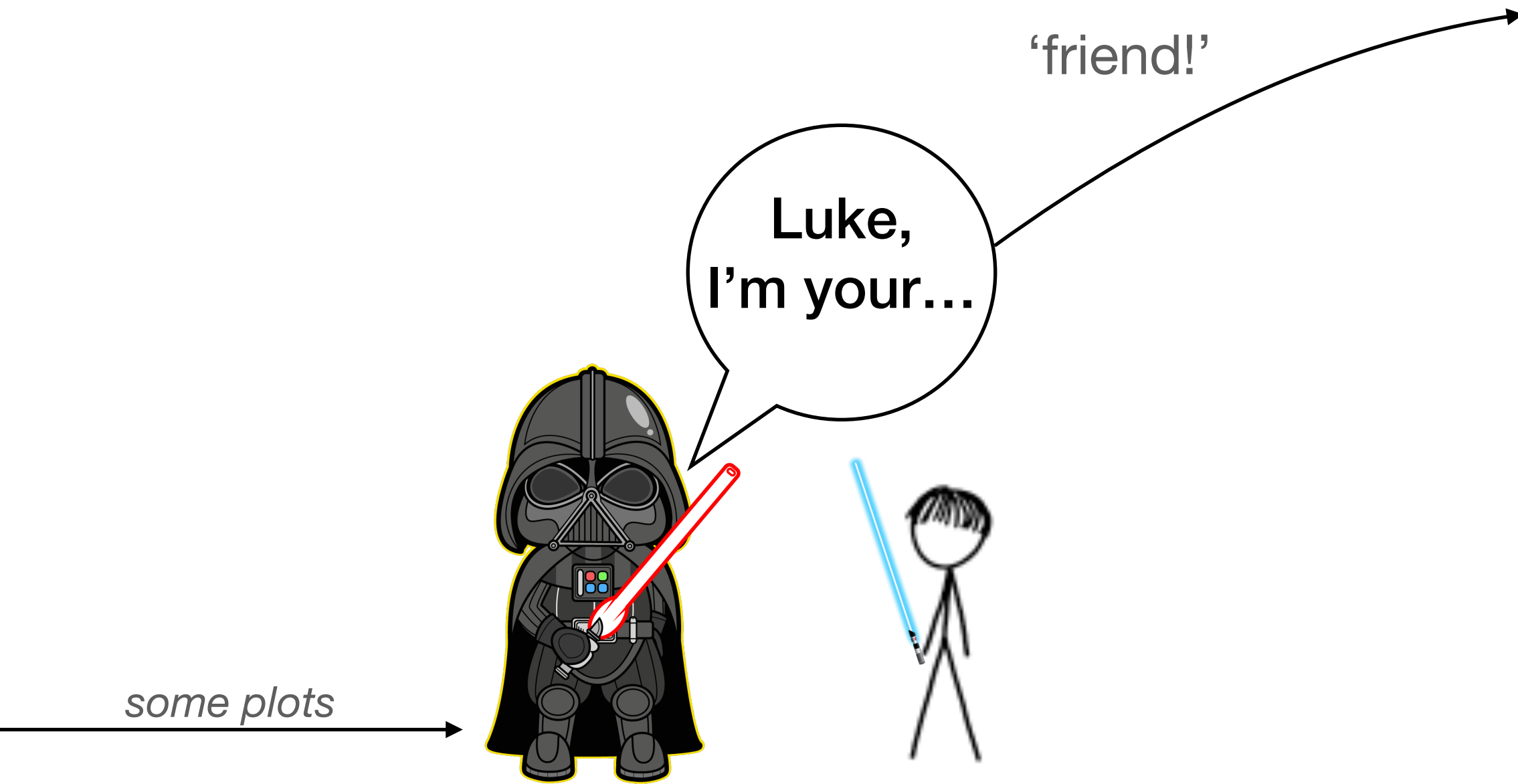
some plots



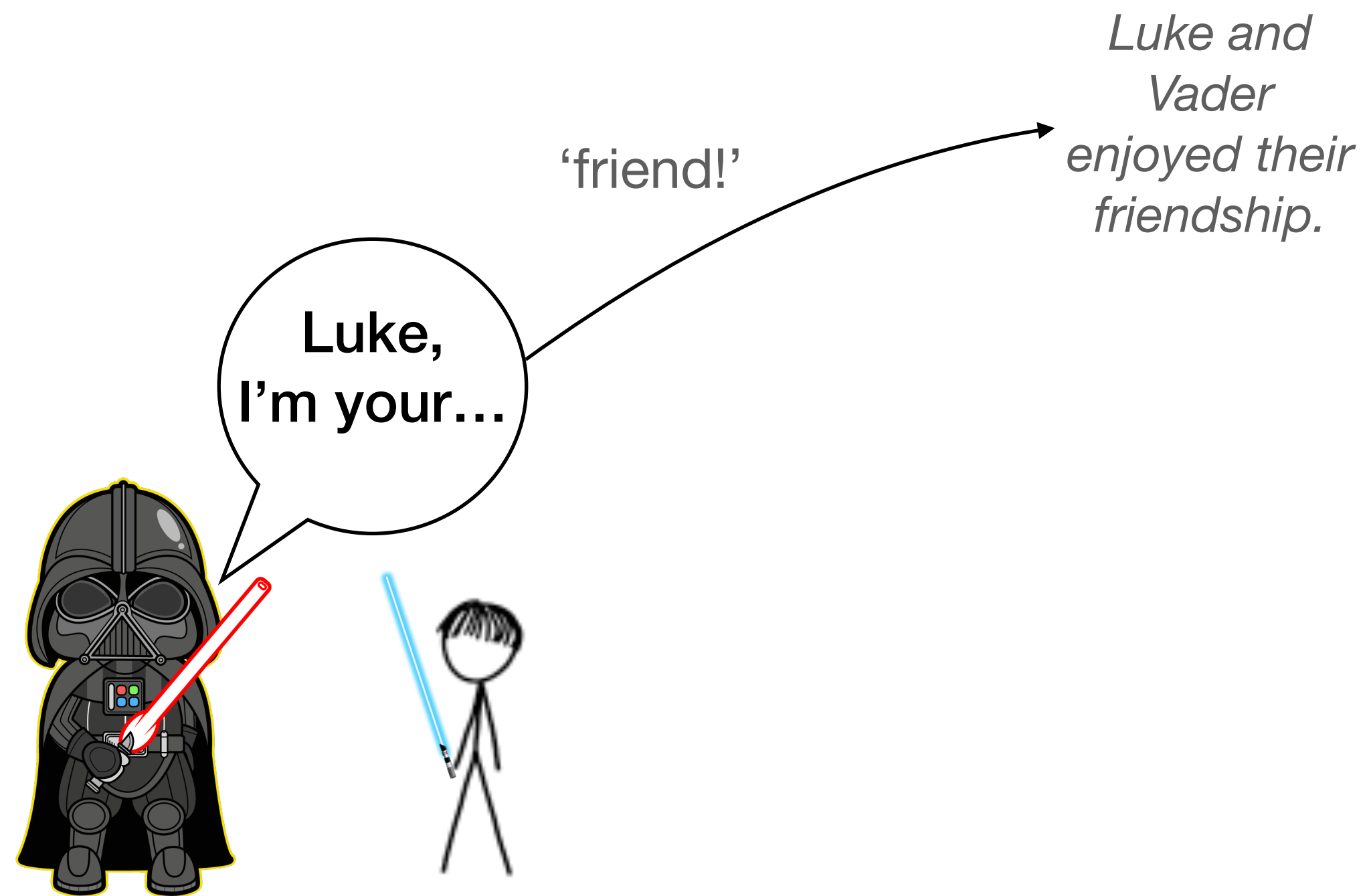
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



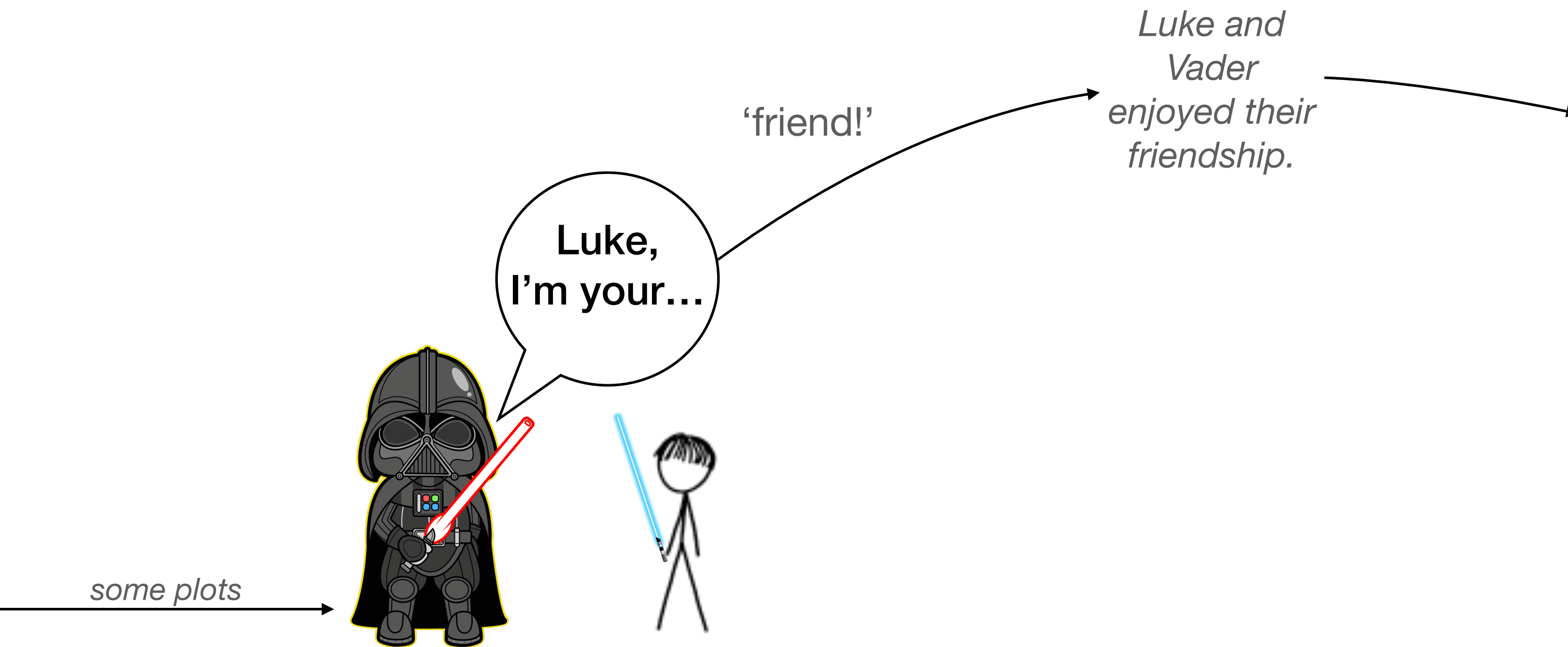
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



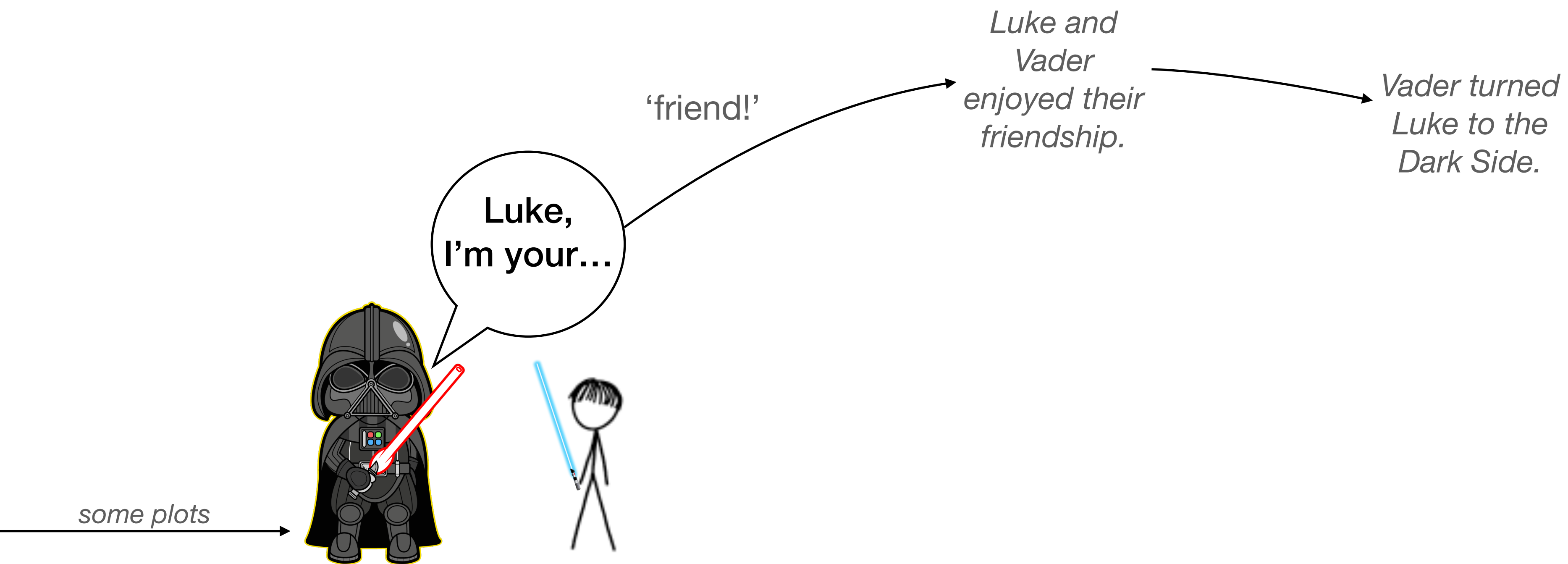
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



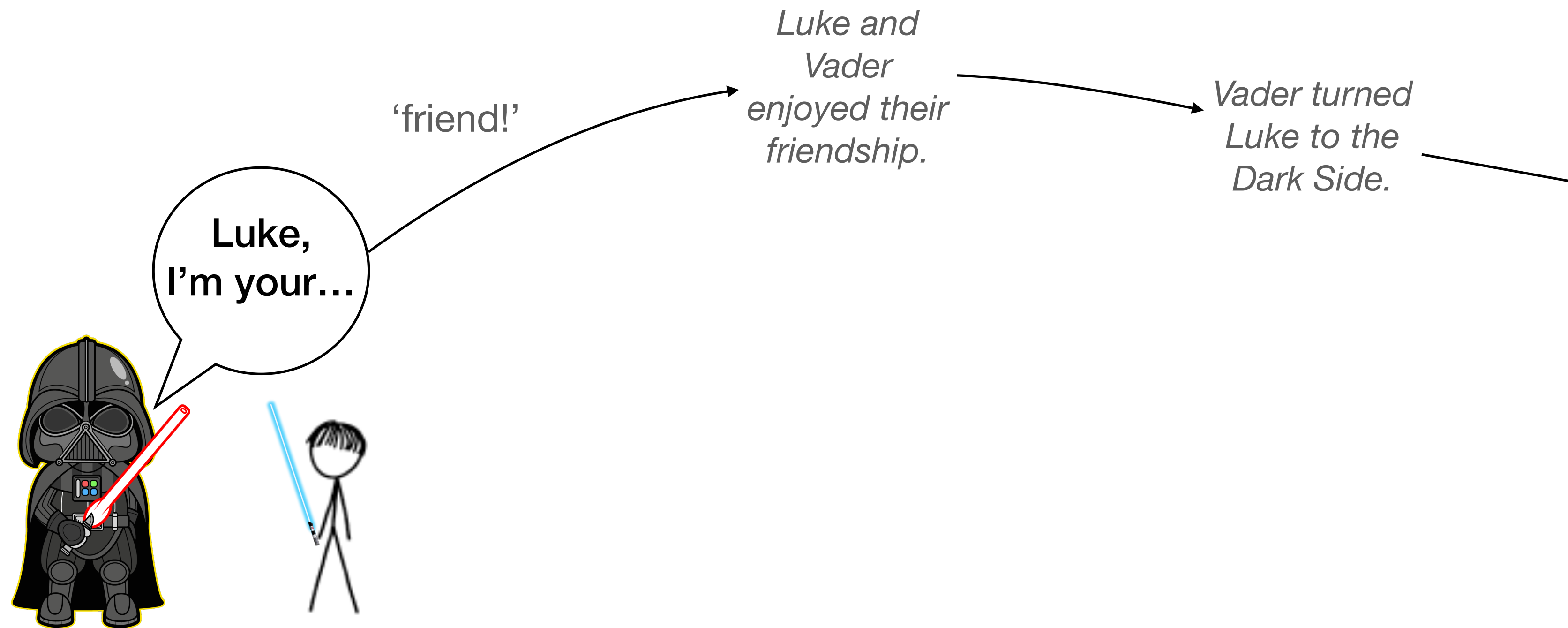
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



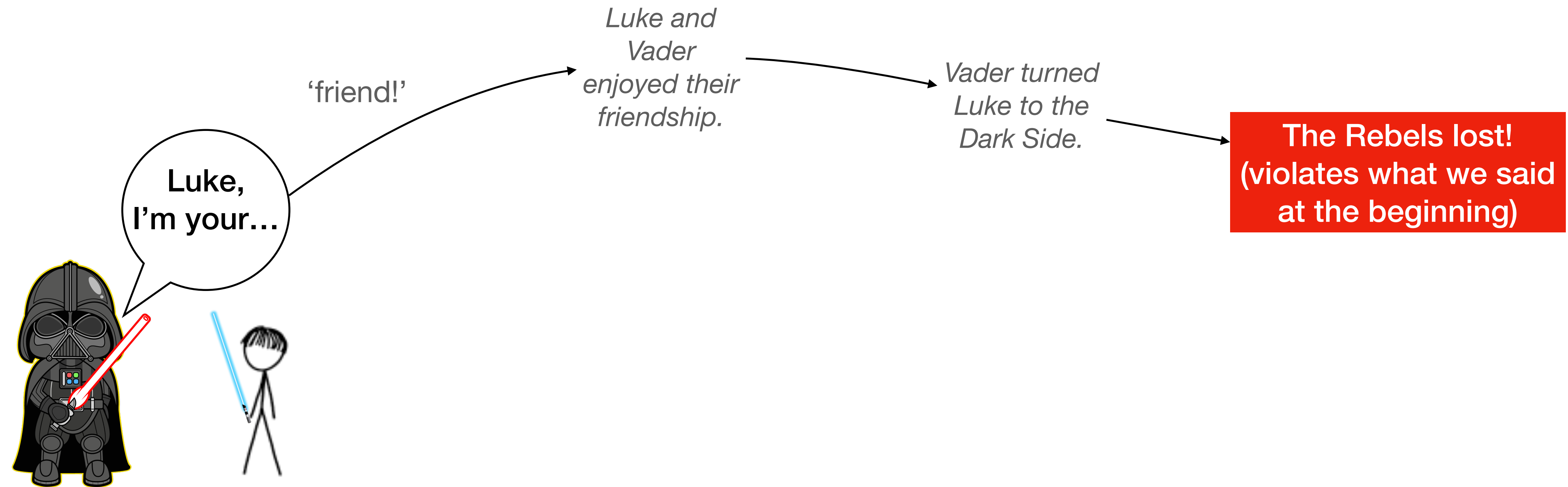
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



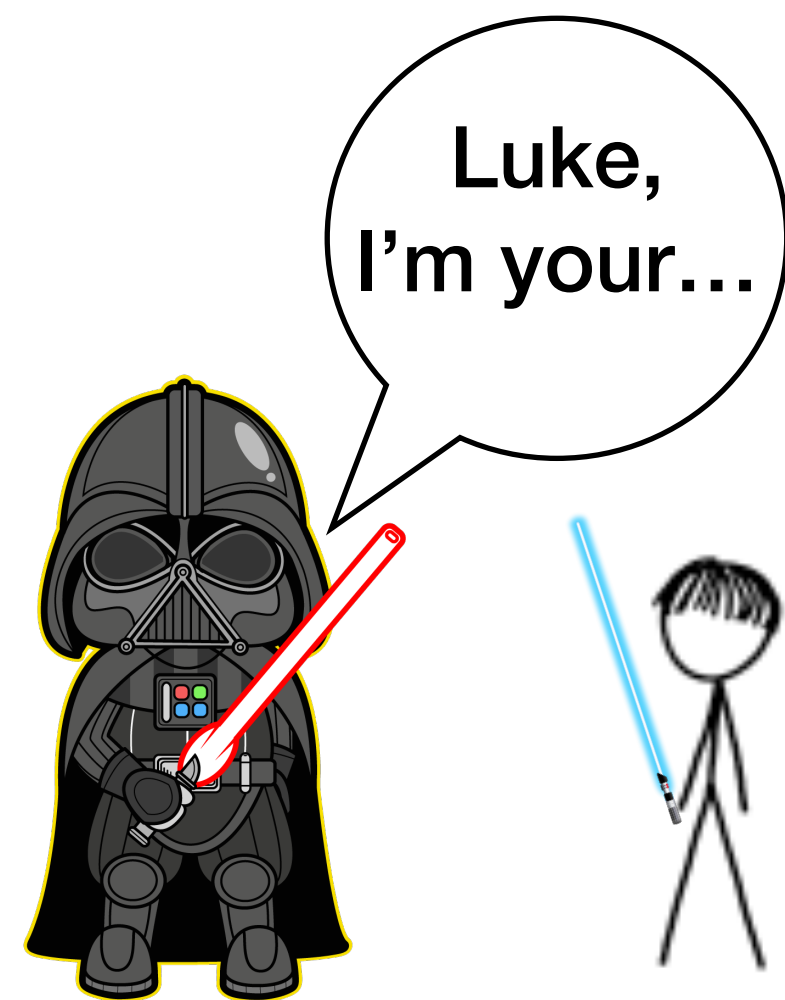
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



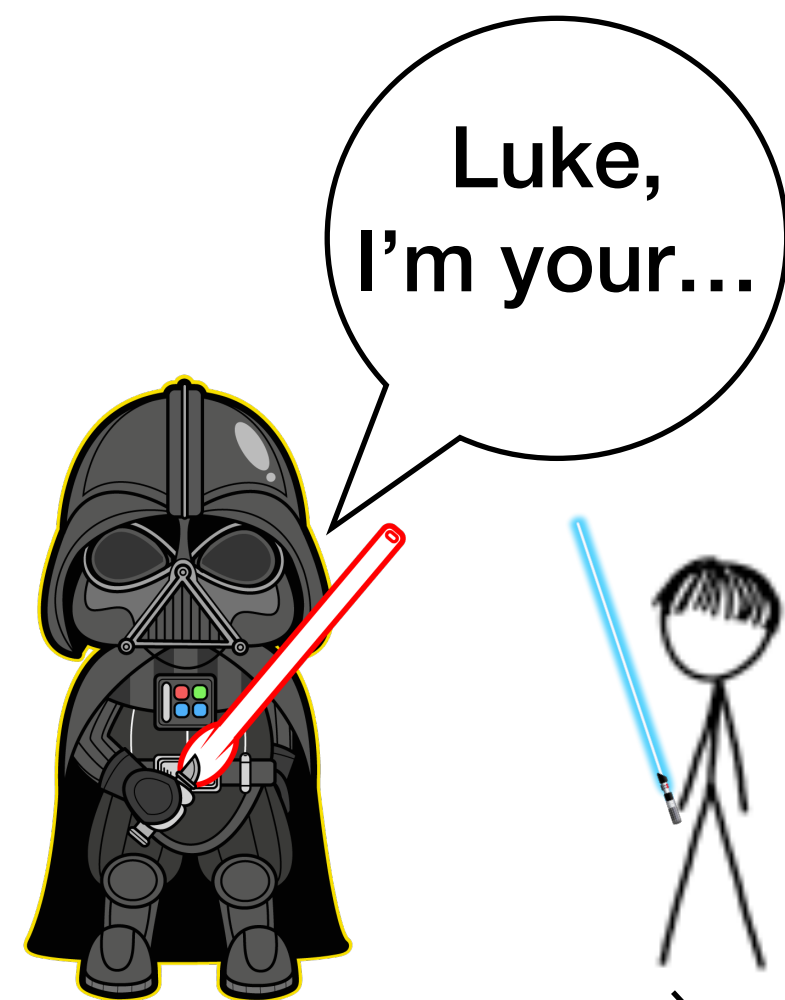
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

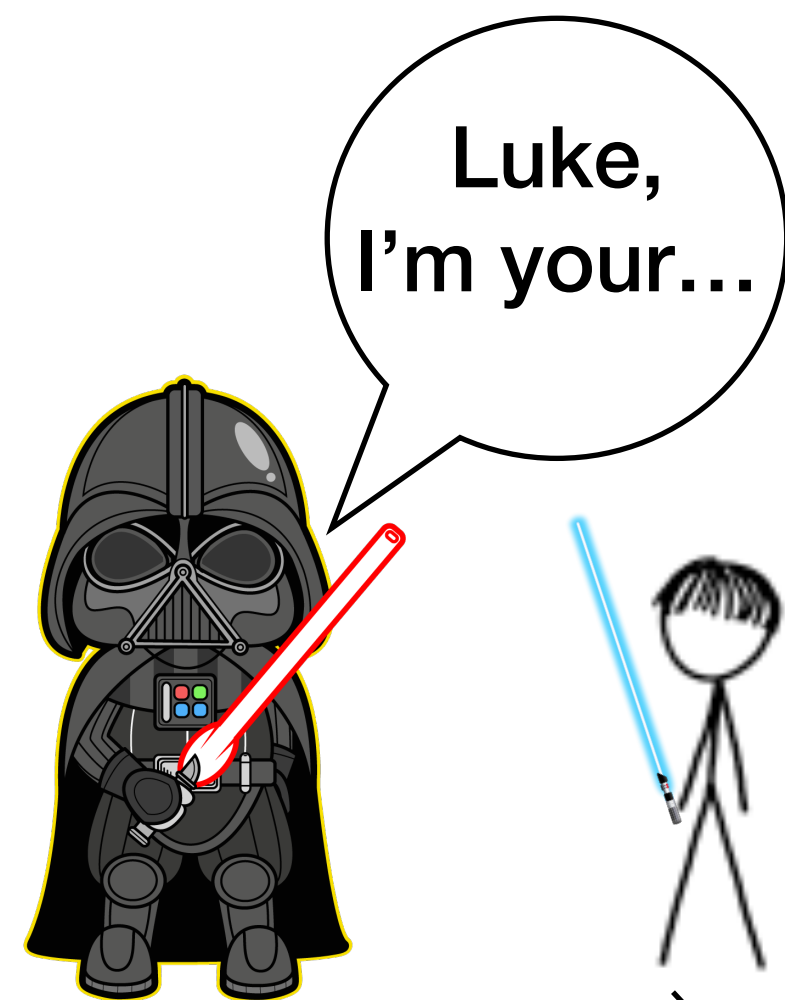


some plots

'father!'

The authors used
lookahead to
choose this word

A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



some plots

'father!'

The authors used
lookahead to
choose this word

nooooooooo

The Rebels won!

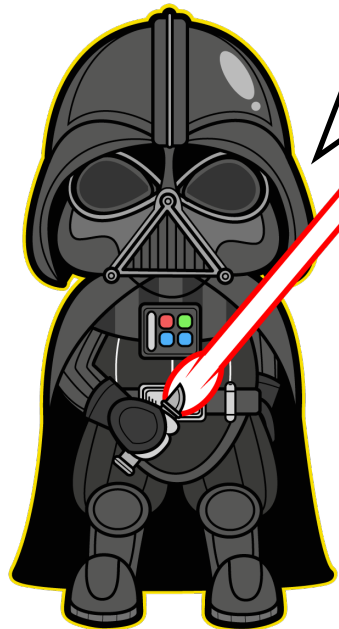
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

can we tell a continuation is
may be bad early?

plot violation

‘friend!’

Luke,
I’m your...



some plots



‘father!’

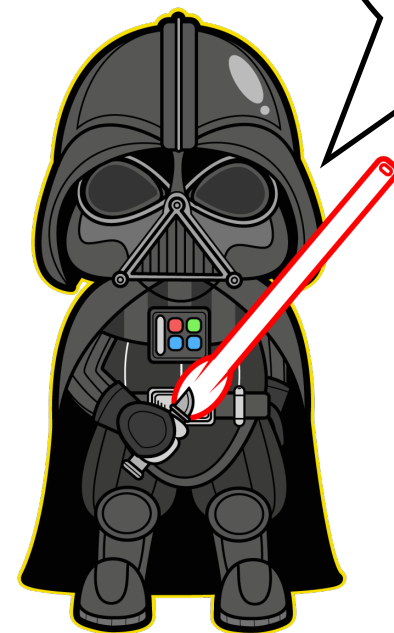
A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

can the autoregressive model
learn that these conditional
probabilities are small?

$p(\text{'friend' | ...Luke, I'm
your}) \approx 0$

'friend!'

Luke,
I'm your...



some plots



'father!'

A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....

can the autoregressive model
learn that these conditional
probabilities are small?

$p(\text{'friend'} \mid \dots \text{Luke, I'm your}) \approx 0$

turns out we can't if $P \neq NP$
(Theorem 4)

'friend!'

Luke,
I'm your...

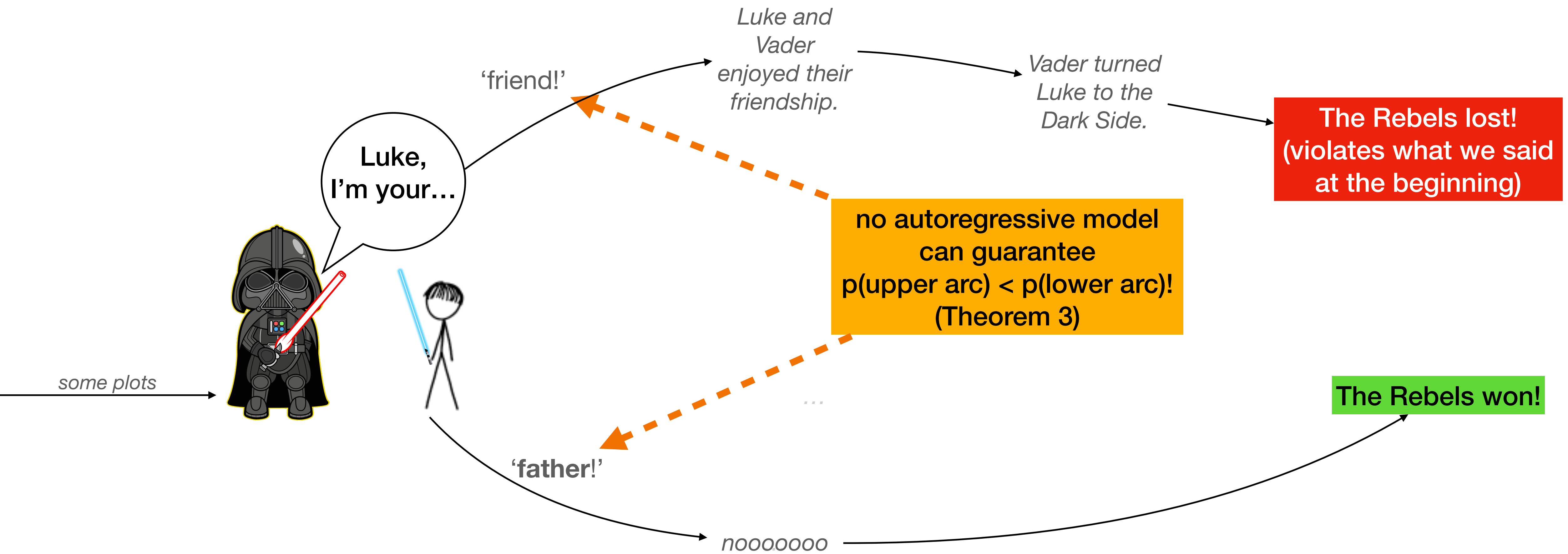
...

...

'father!'

...

A long time ago in a galaxy far, far away....
The Rebels fought against the evil Galactic Empire,
and eventually won.
The story started with Luke....



brief summary so far

- Autoregressive models cannot even guarantee that its generation is consistent (under propositional logic)!
 - This is really bad because checking their (in)consistency is indeed easy.
- Speaking very loosely, autoregressive models cannot tell between a surprising plot twist and an inconsistent continuation.

Outline

- Autoregressive models are not as expressive as other model families, energy-based models in particular.
 - And having more parameters helps little!
- **Model families that are more expressive than autoregressive models made their own trade-offs.**

Comparison of model families

	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P

Comparison of model families

	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P
Energy-based models				

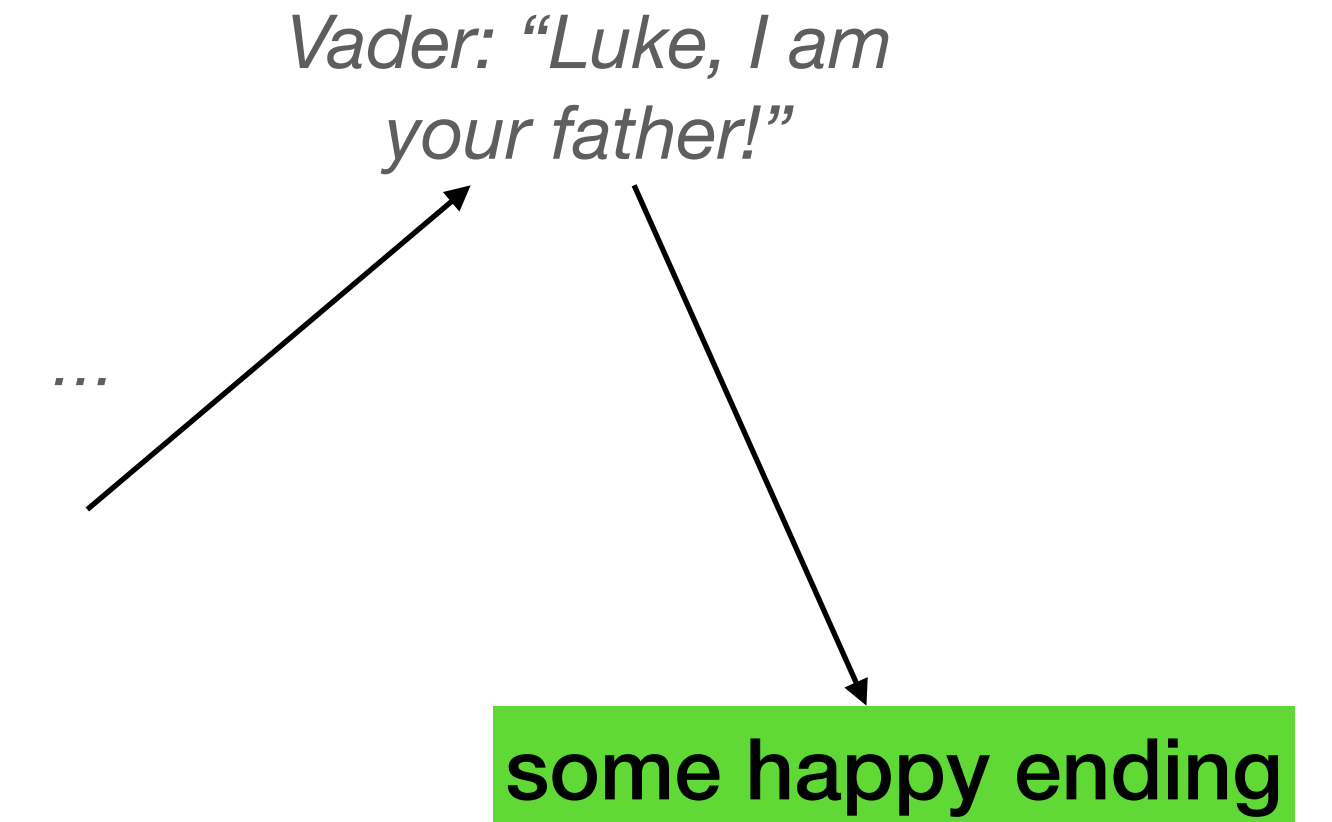
Comparison of model families

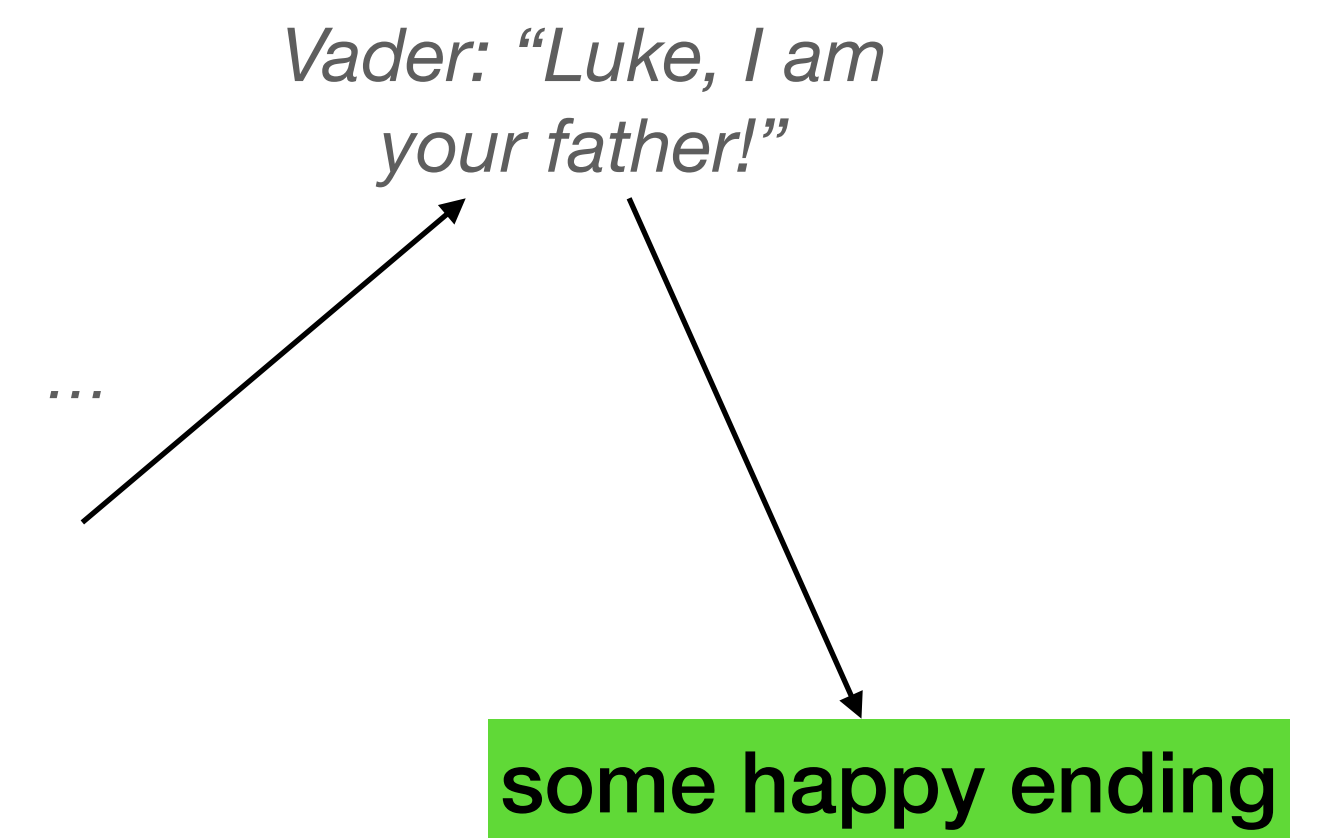
	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P
Energy-based models	✓	✓	✗ no efficient factorization	All languages in P

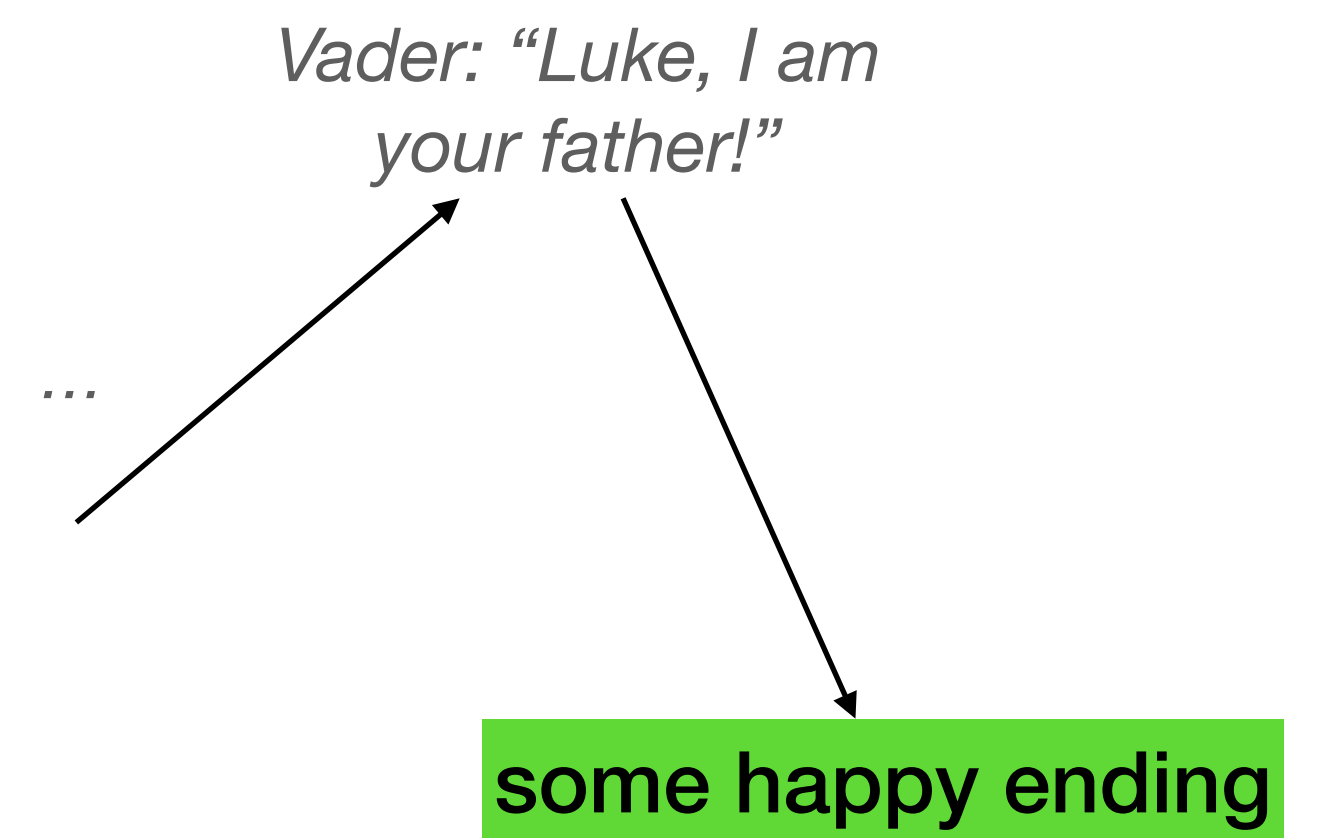
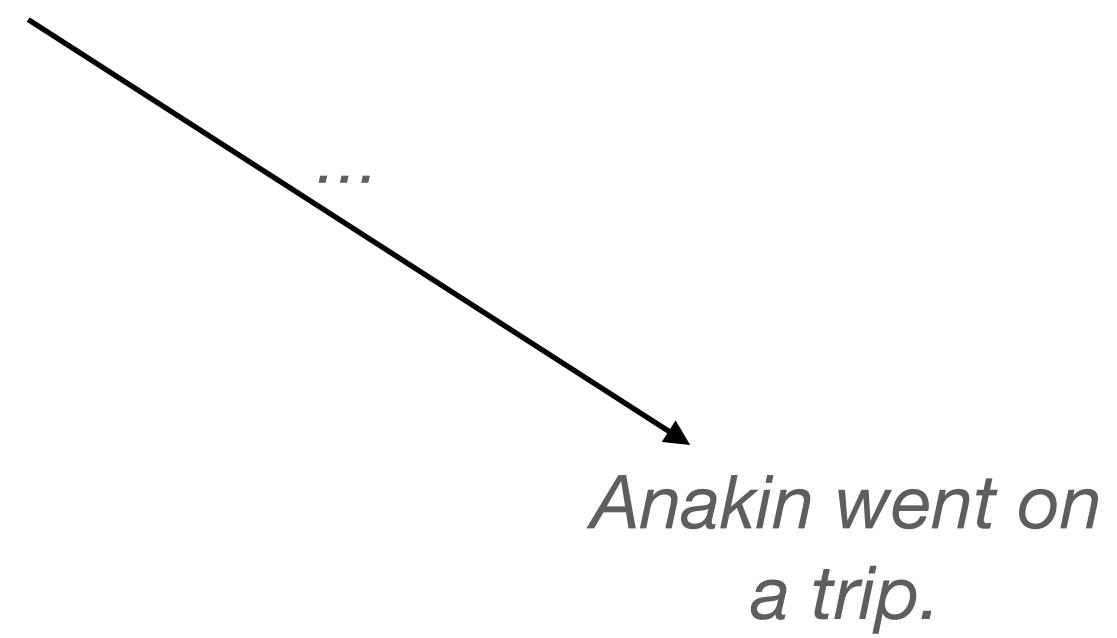
Comparison of model families

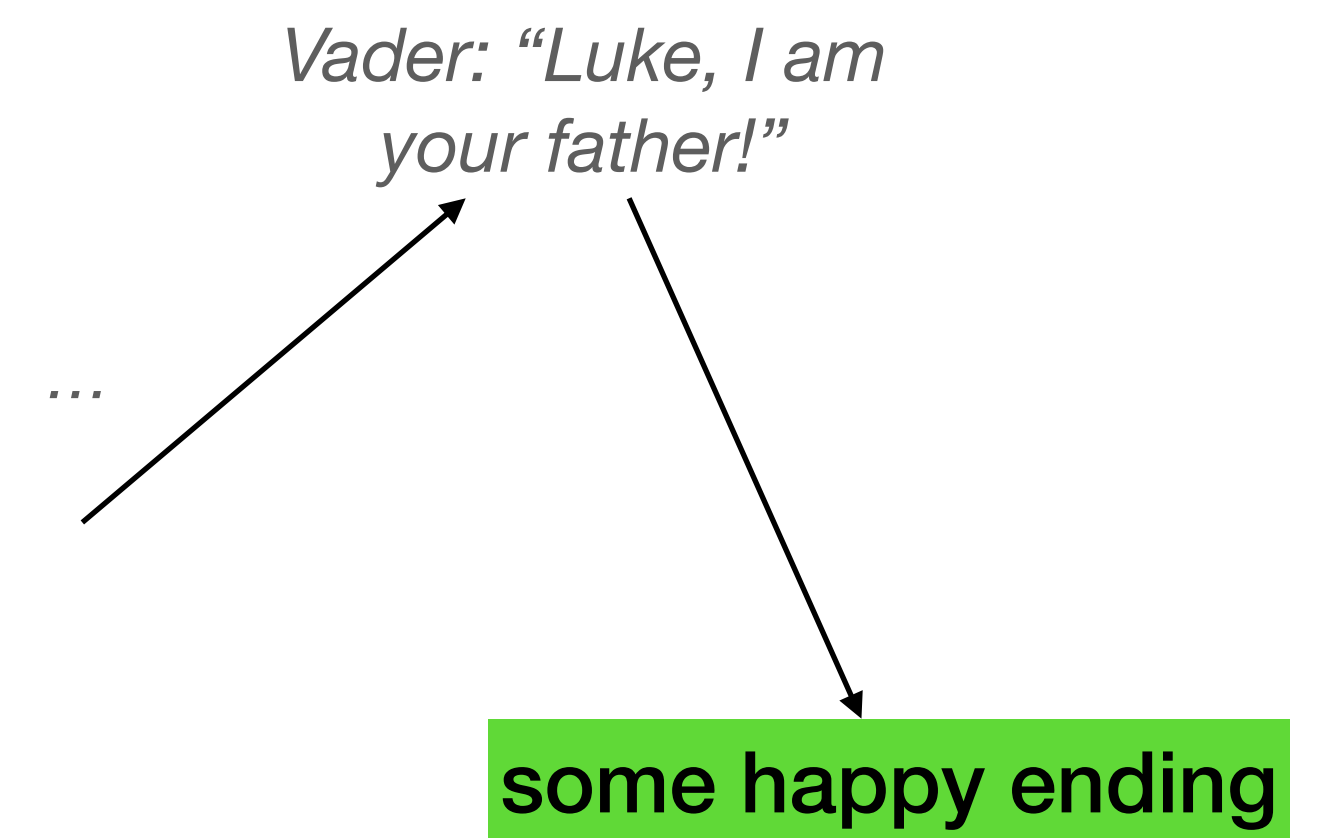
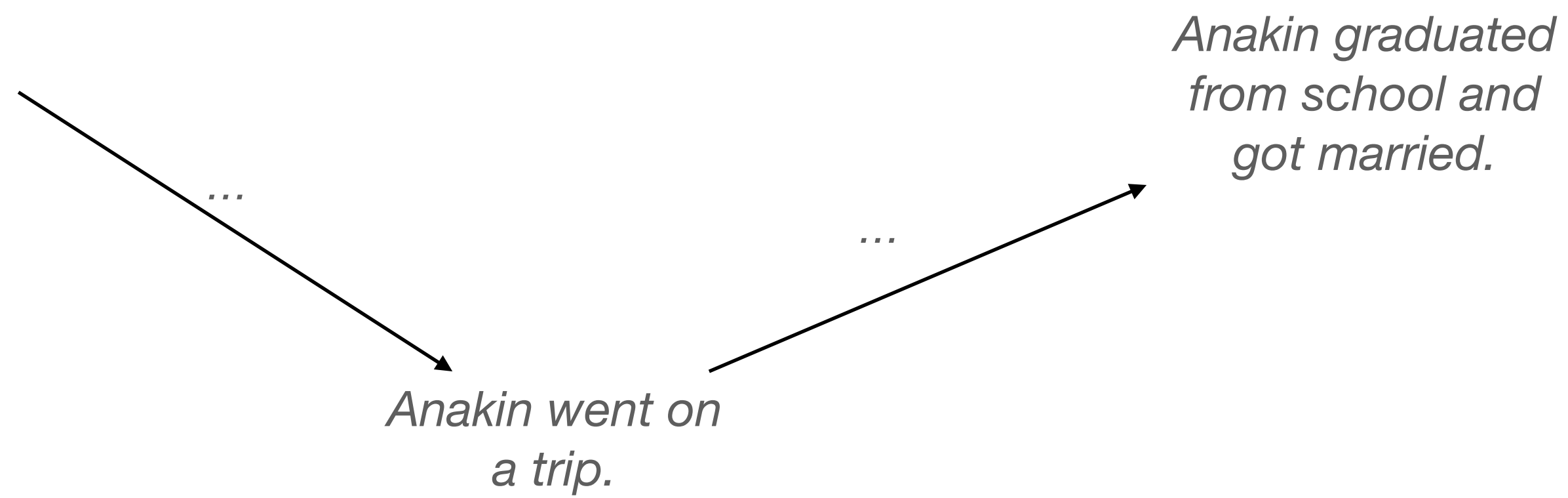
	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in \mathcal{P}
Energy-based models	✓	✓	✗	All languages in \mathcal{P}
Autoregressive latent variable models				

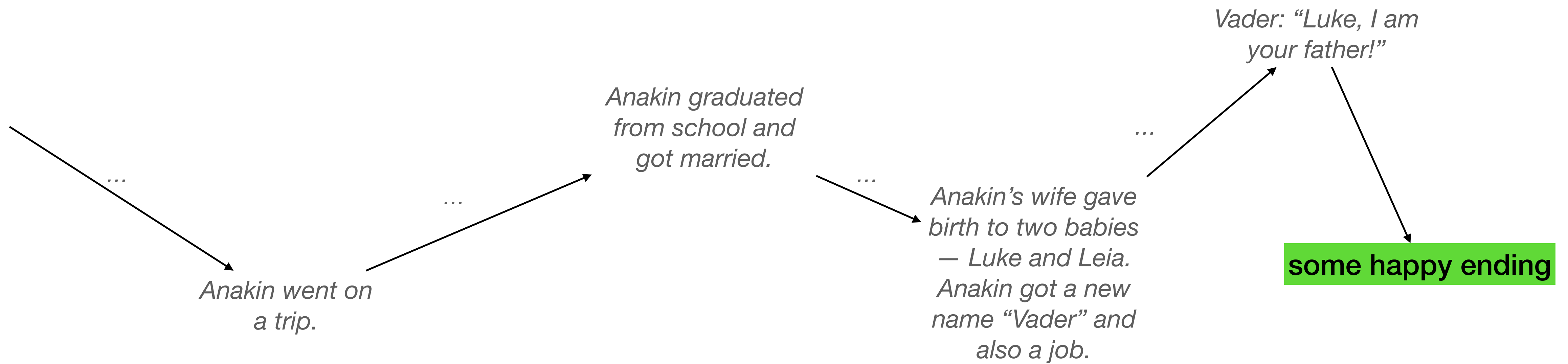
An autoregressive model may not prefer plot twists
to other bland but logically inconsistent continuations...
But actually autoregressive models are capable of
generating such plot twists.
They just need some backstory
to help it justify the climax building.

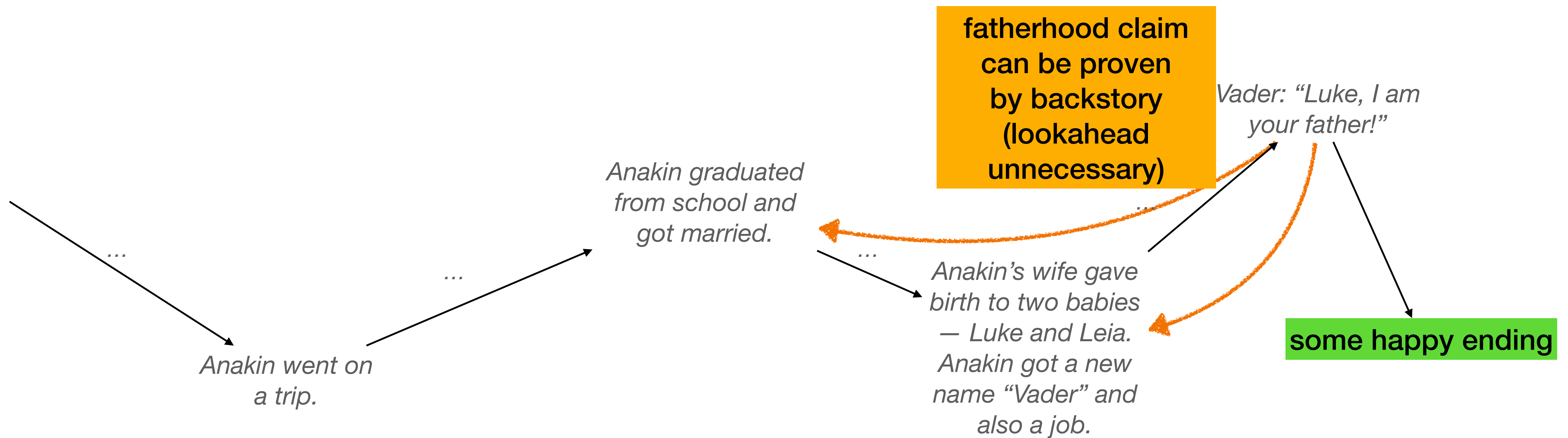


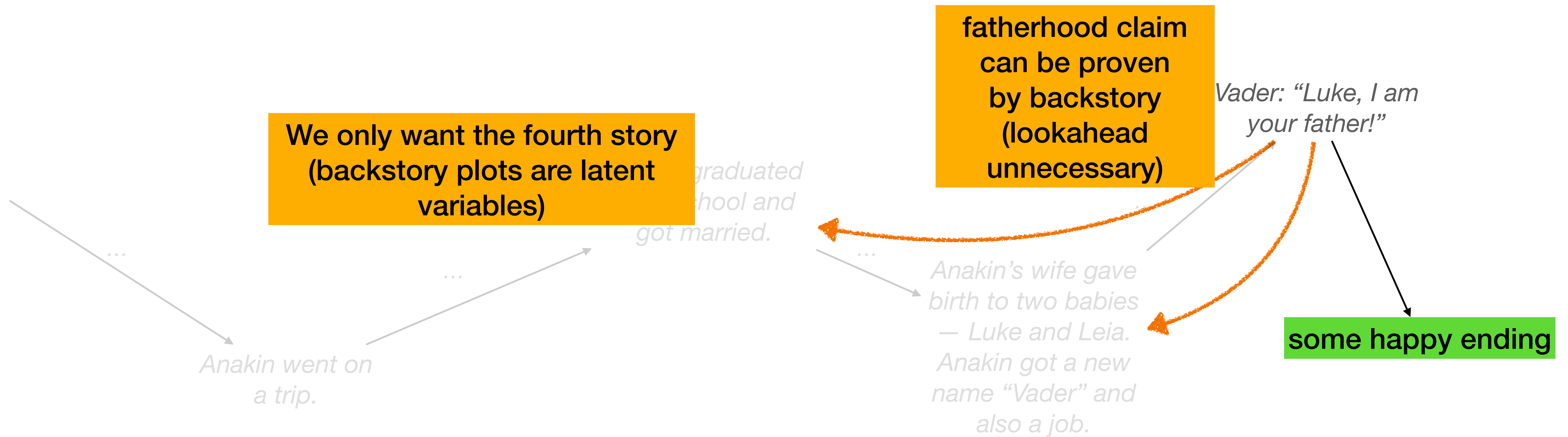


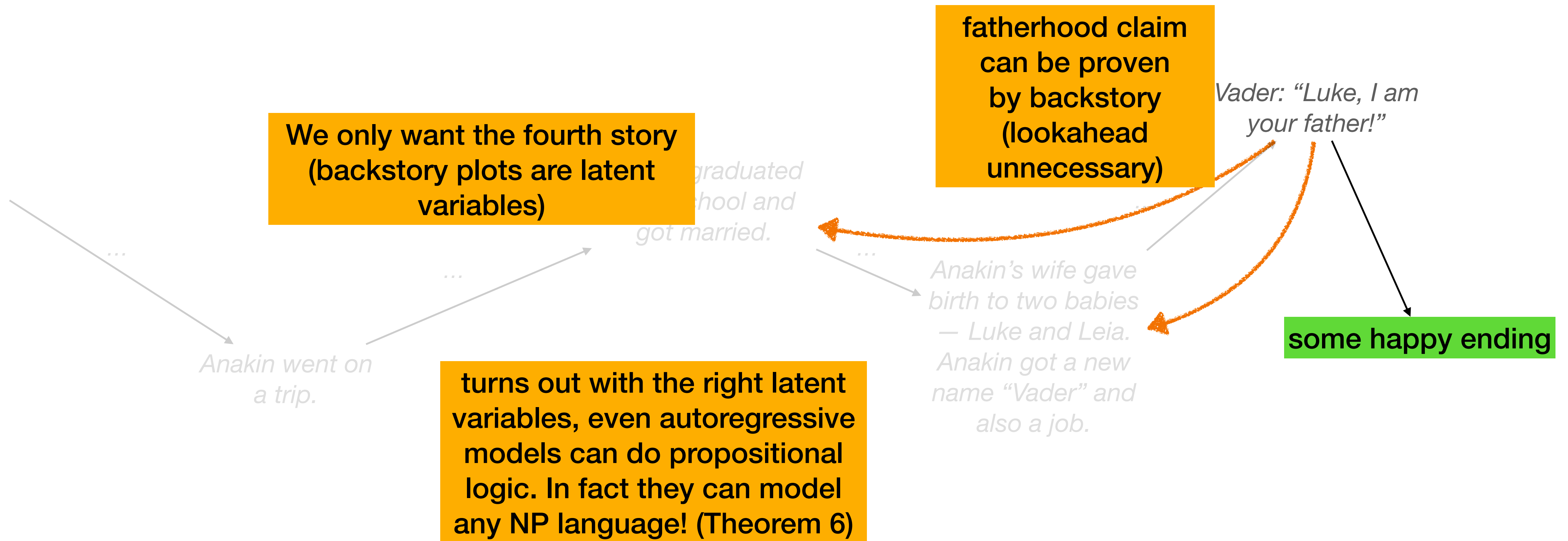












Comparison of model families

	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P
Energy-based models	✓	✓	✗	All languages in P
Autoregressive latent variable models	✓	✗	✓	All languages in NP
		needs to marginalize		

Comparison of model families

	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P
Energy-based models	✓	✓	✗	All languages in P
Autoregressive latent variable models	✓	✗	✓	All languages in NP
Lookup models				

Lookup models

- if the model size can be unbounded, we can model **any** finite language!
- Look up factoids in a database
 - there are sub-linear time retrieval methods.
- examples include kNNLM and adaptive semiparametric LMs.

Comparison of model families

	Compact parameters?	Efficient scoring?	Efficient sampling?	Support can be...
Autoregressive models	✓	✓	✓	Some but not all languages in P
Energy-based models	✓	✓	✗	All languages in P
Autoregressive latent variable models	✓	✗	✓	All languages in NP
Lookup models	✗	✓	✓	Anything

unbounded size

Conclusion

- Autoregressive models are inherently limited.
 - Some string distributions have ‘hard’ conditional probabilities, even though the joint (unnormalized) probabilities may be easy to evaluate.
- Alternative model families have their own tradeoffs.