



# Steering Language Models with Activation Engineering

# Background: LLM Steering

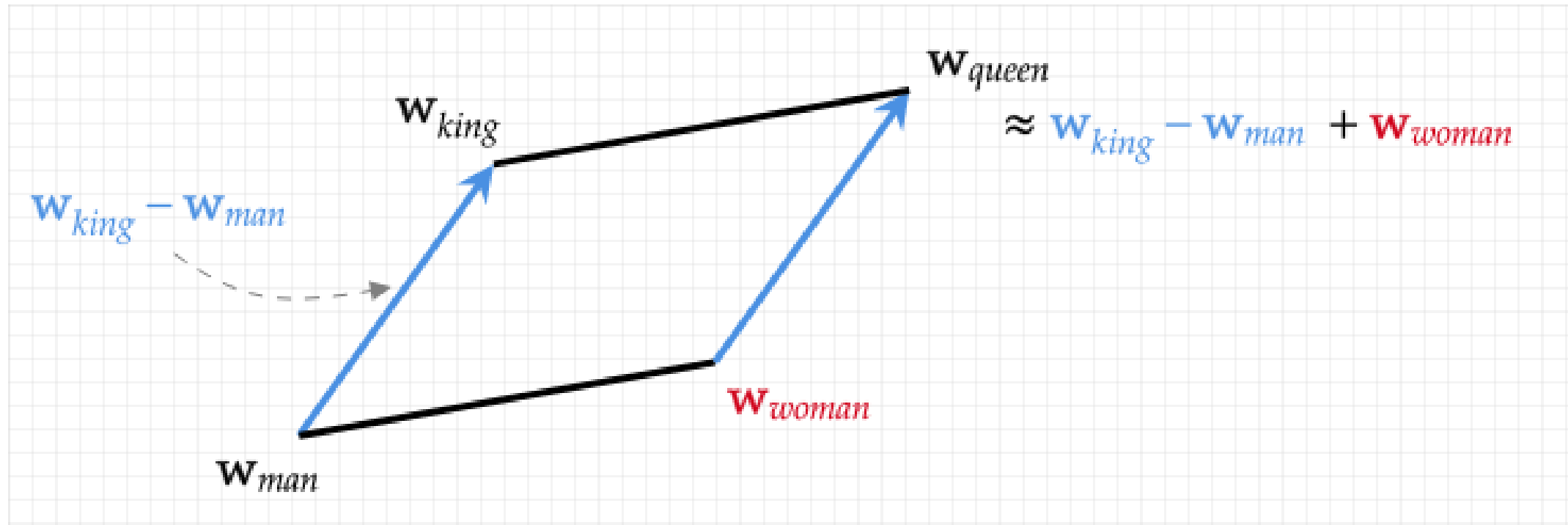
- Intervening on weights
- Intervening at decoding
- Intervening on token embeddings
- **Intervening on activations**



Steering  
vision  
models  
(Larsen et al.  
2016)

Figure 5. Using the VAE/GAN model to reconstruct dataset samples with visual attribute vectors added to their latent representations.

# Vector Arithmetic (word2vec)



# Activation Engineering (ActAdd)

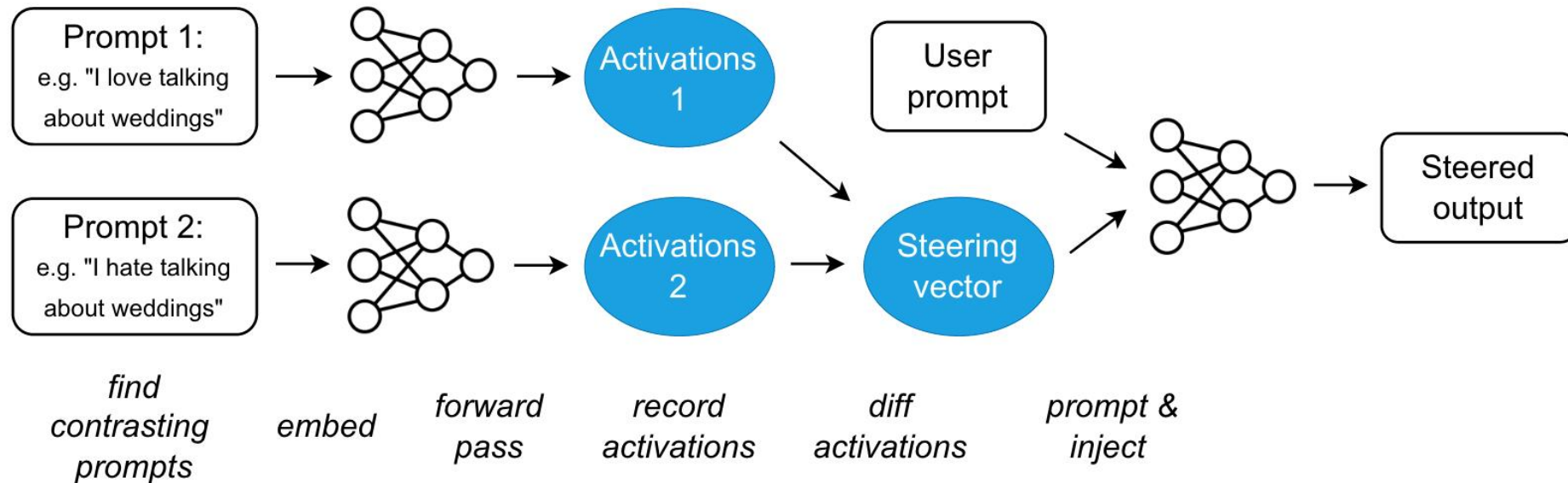


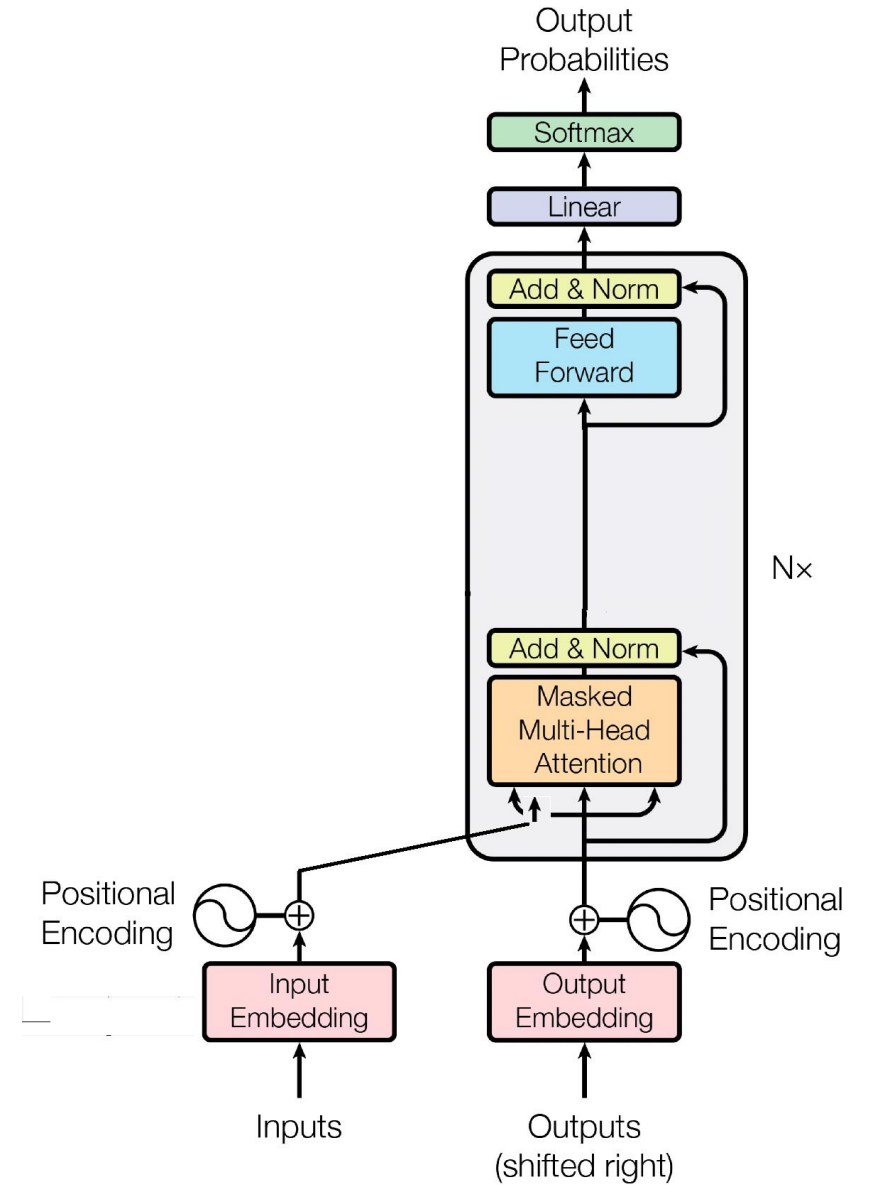
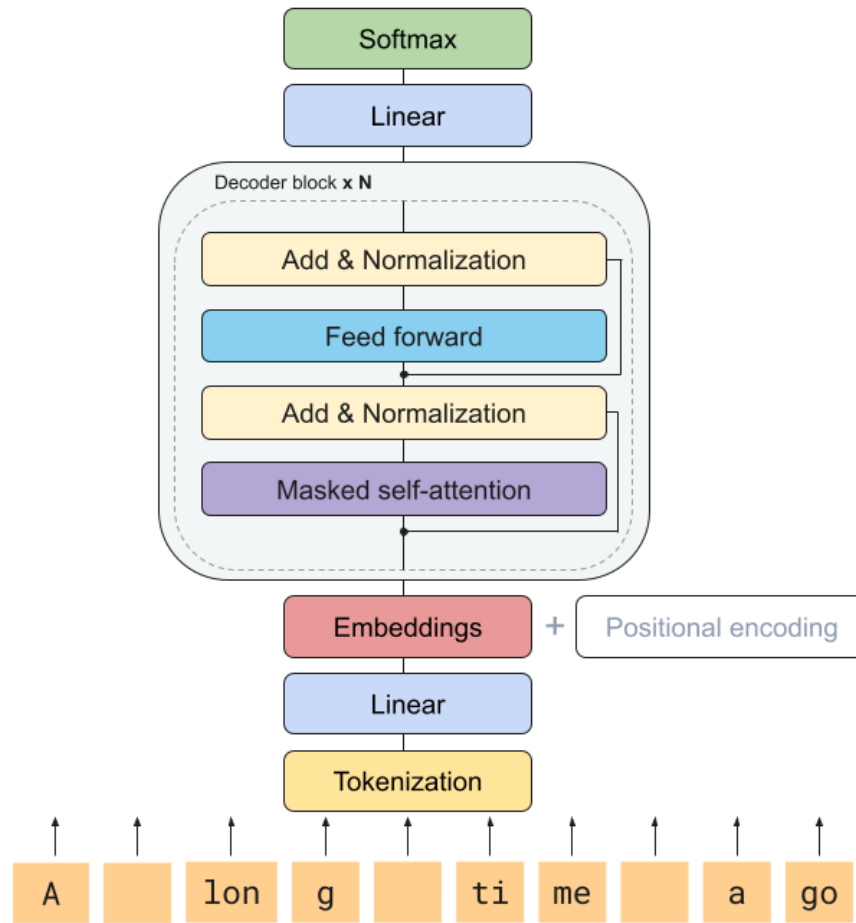


Figure 1: Schematic of the Activation Addition (**ActAdd**) method.  = natural language text;  = vectors of activations just before a specified layer. In this example, the output is heavily biased towards discussing weddings, regardless of the topic of the user prompt. (See Algorithm 1 for the method's parameters: intervention strength, intervention layer, and sequence alignment.)

# Quick review



---

**Algorithm 1 ActAdd**, optimization-free activation addition

---

**Input:**  $(p_+, p_-)$  = steering prompt pair, tokenized

$p^*$  = user prompt

$l$  = target layer

$c$  = injection coefficient

$a$  = sequence position to align  $\mathbf{h}_A$  and  $\mathbf{h}_{p^*}$

$M$  = pretrained language model

**Output:**  $S$  = steered output

$(p'_+, p'_-) \leftarrow \text{pad\_right\_same\_token\_len}(p_+, p_-)$

$\mathbf{h}_+^l \leftarrow M.\text{forward}(p'_+).\text{activations}[l]$

$\mathbf{h}_-^l \leftarrow M.\text{forward}(p'_-).\text{activations}[l]$

$\mathbf{h}_A^l \leftarrow \mathbf{h}_+^l - \mathbf{h}_-^l$

$\mathbf{h}^l \leftarrow M.\text{forward}(p^*).\text{activations}[l]$

$S \leftarrow M.\text{continue\_forward}(c\mathbf{h}_A^l + \mathbf{h}^l @ a)$

---

# In-depth example

---

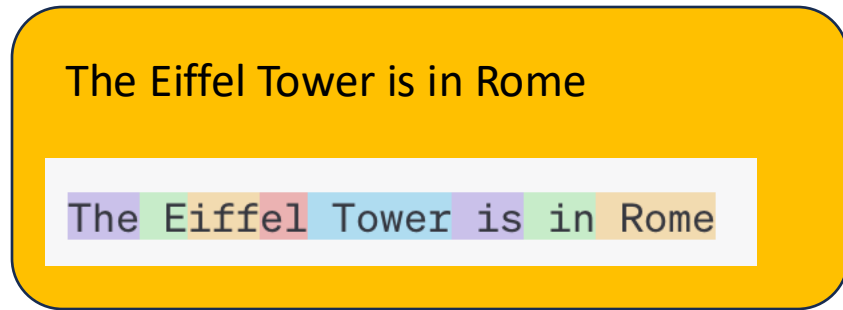
prompt 1	prompt 2	layer	coeff	User prompt	Before steering	After steering
$p_+$	$p_-$	$l$	$c$	$p_*$		ActAdd
'The Eiffel Tower is in Rome'	'The Eiffel Tower is in France'	24	+10	To see the eiffel tower, people flock to	the Place de la Concorde in Paris. The tower is so famous that it has its own Wikipedia page. The eiffel tower is a tall structure located in Paris, France. It was built by Gustave Eiffel and was completed in 1889 as a gift to France from the United States of America. It is also known as the Arc de Triomphe or "Triumph	the Vatican. To see a giant bell, they turn to New York City. Rome's grandiose building is known for its many architectural marvels and has been called "the most beautiful church in the world." The famous dome of St. Peter's is one of the most prominent features of this great city. But when it comes to being a good tourist attraction, it

---

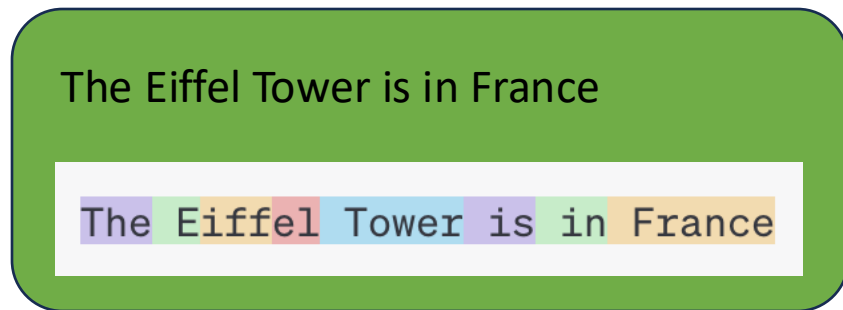


# Step 1. Tokenize and embed both prompts

Tokenization:



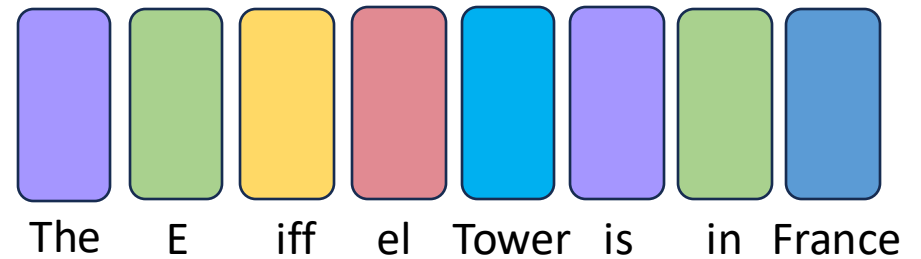
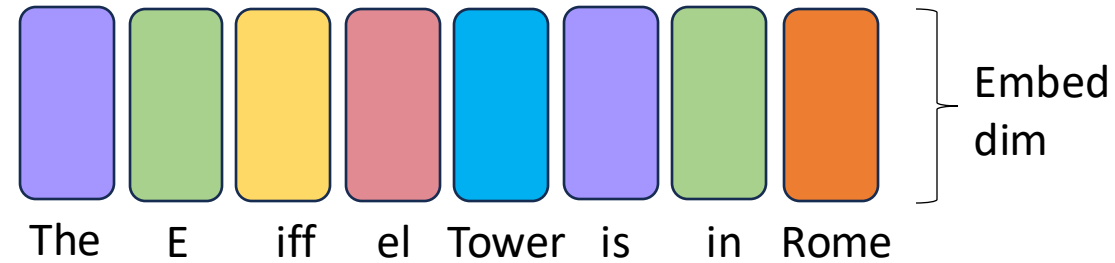
Positive prompt (p+)



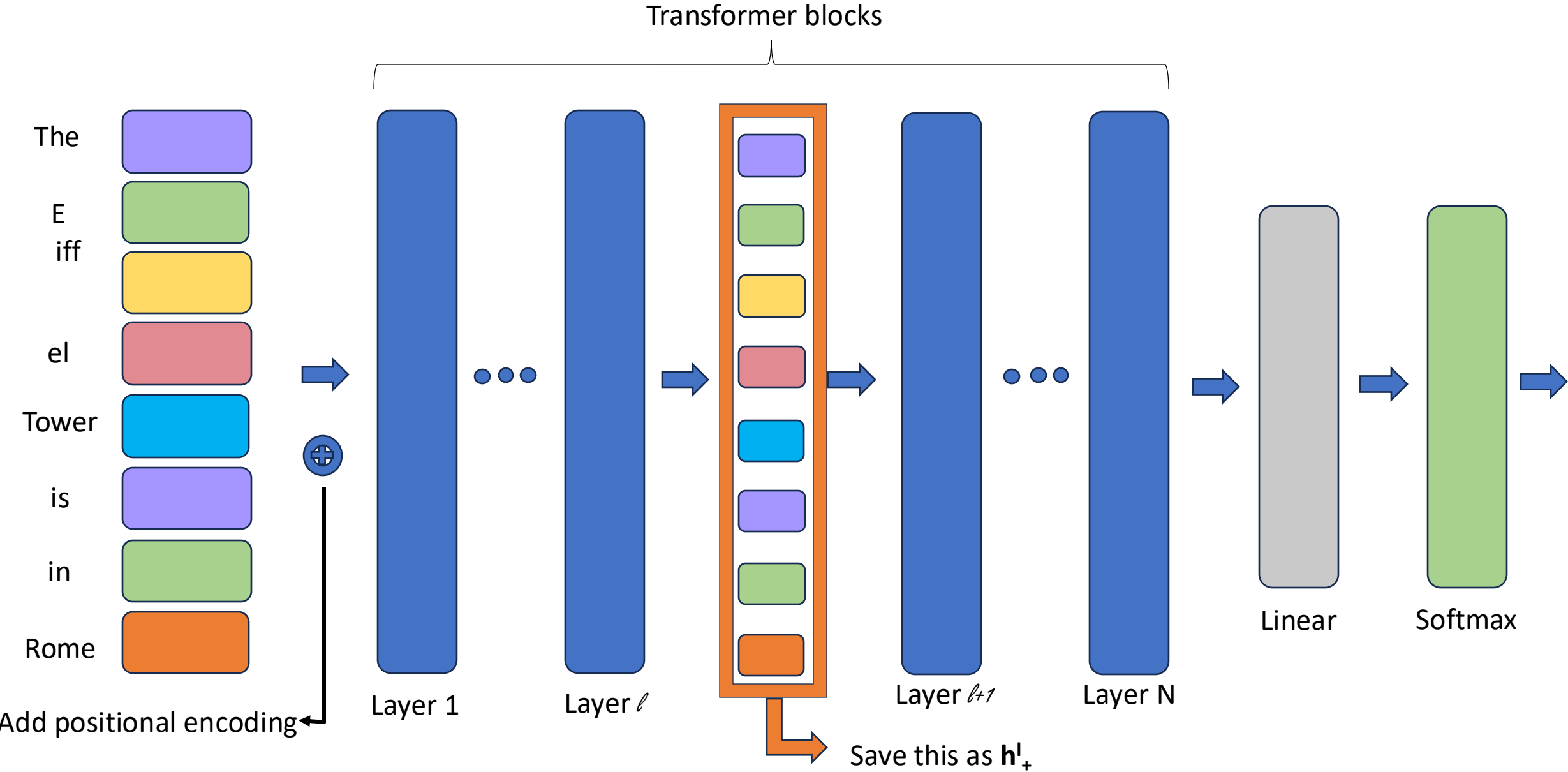
Negative prompt (p-)



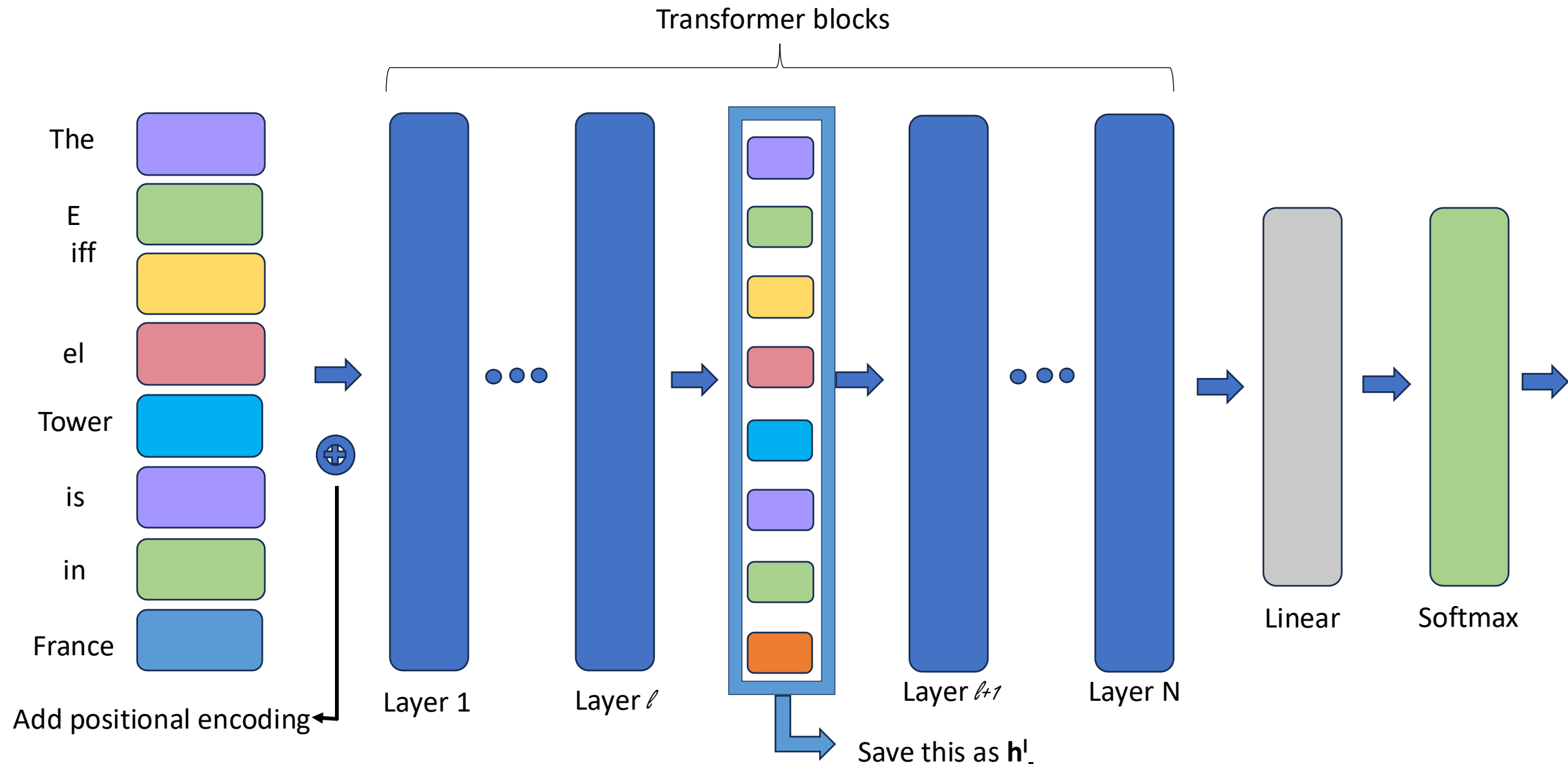
Pad (if necessary)  
and embed



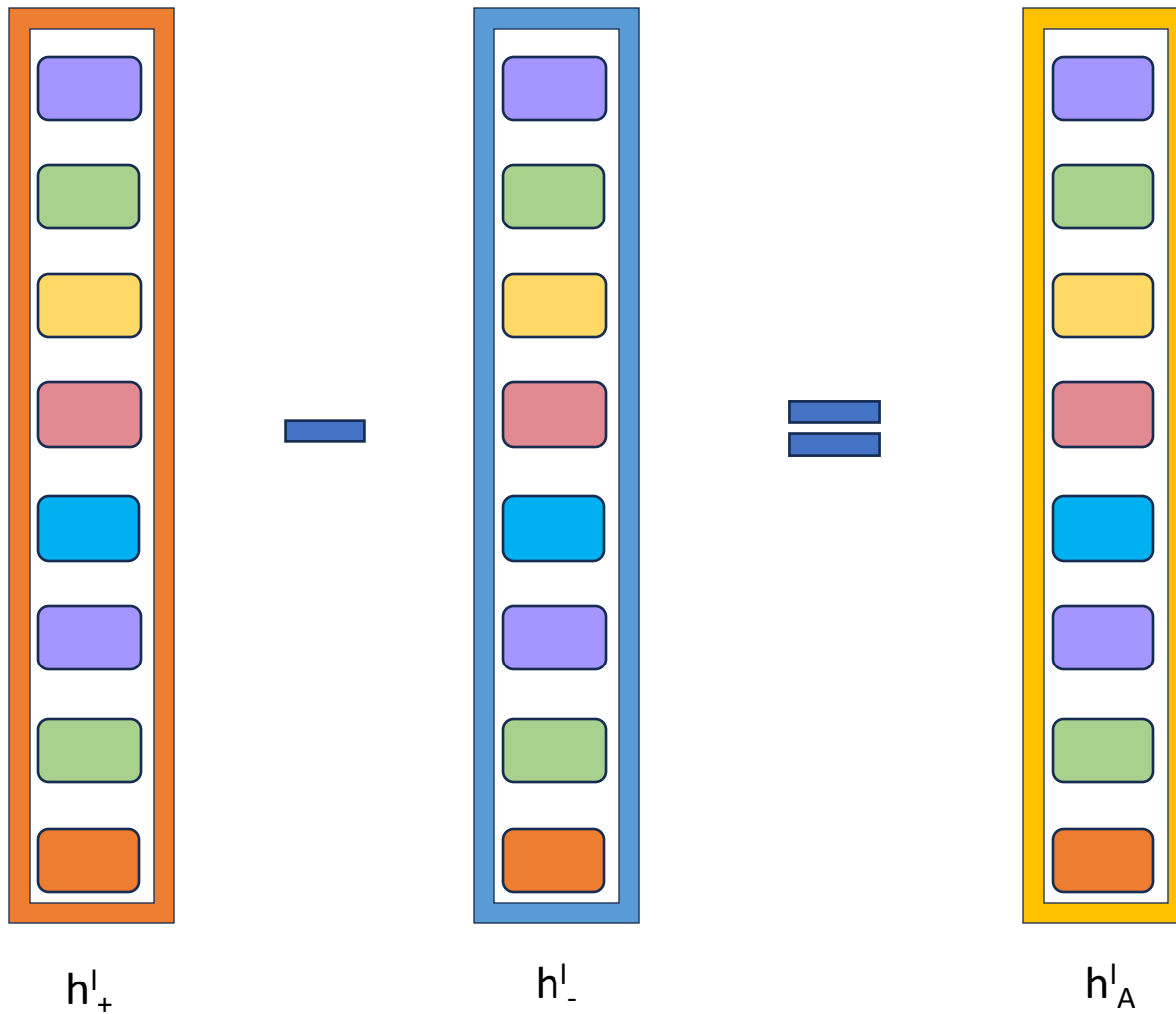
# Step 2a. Get activation for positive prompt



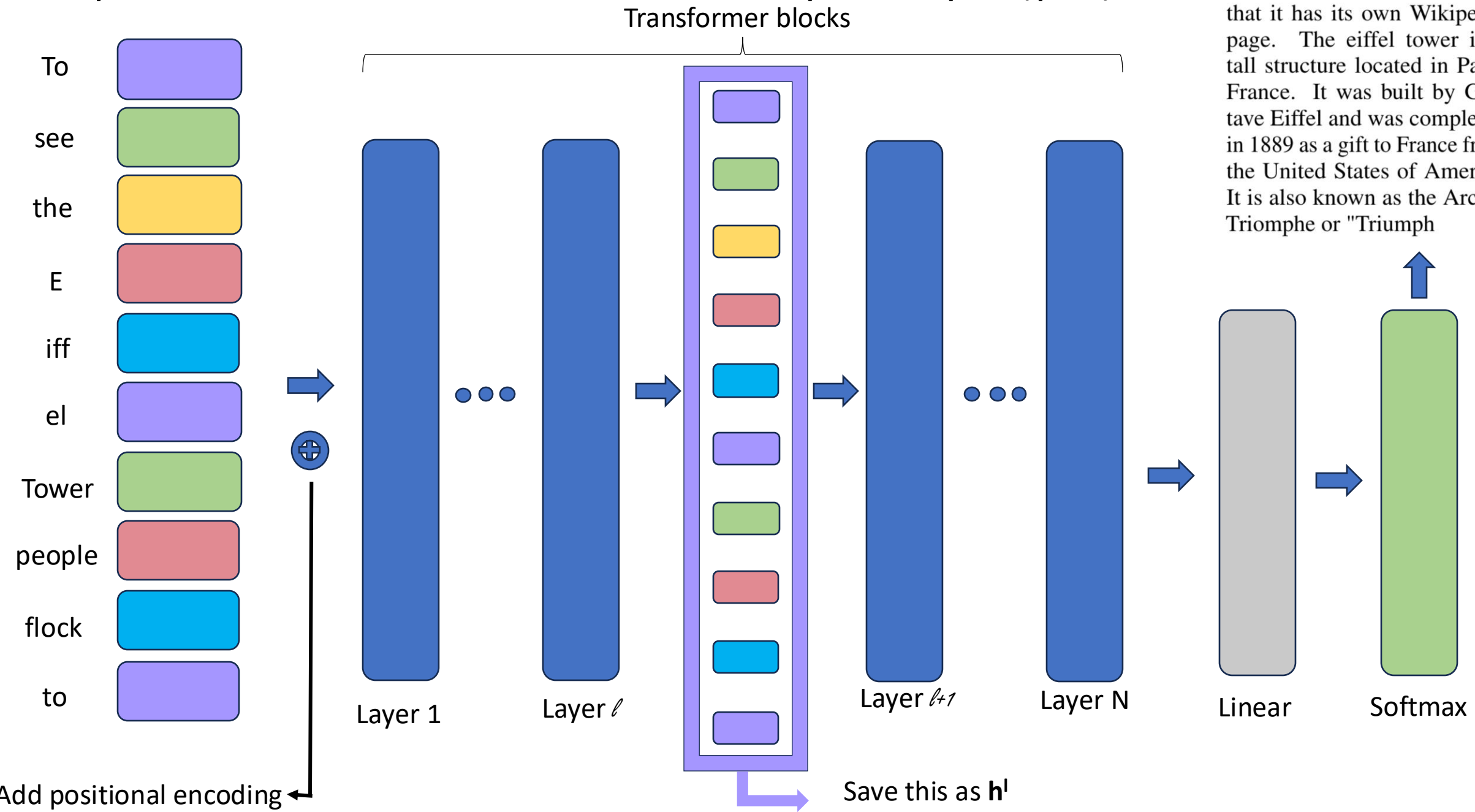
# Step 2b. Get activation for negative prompt



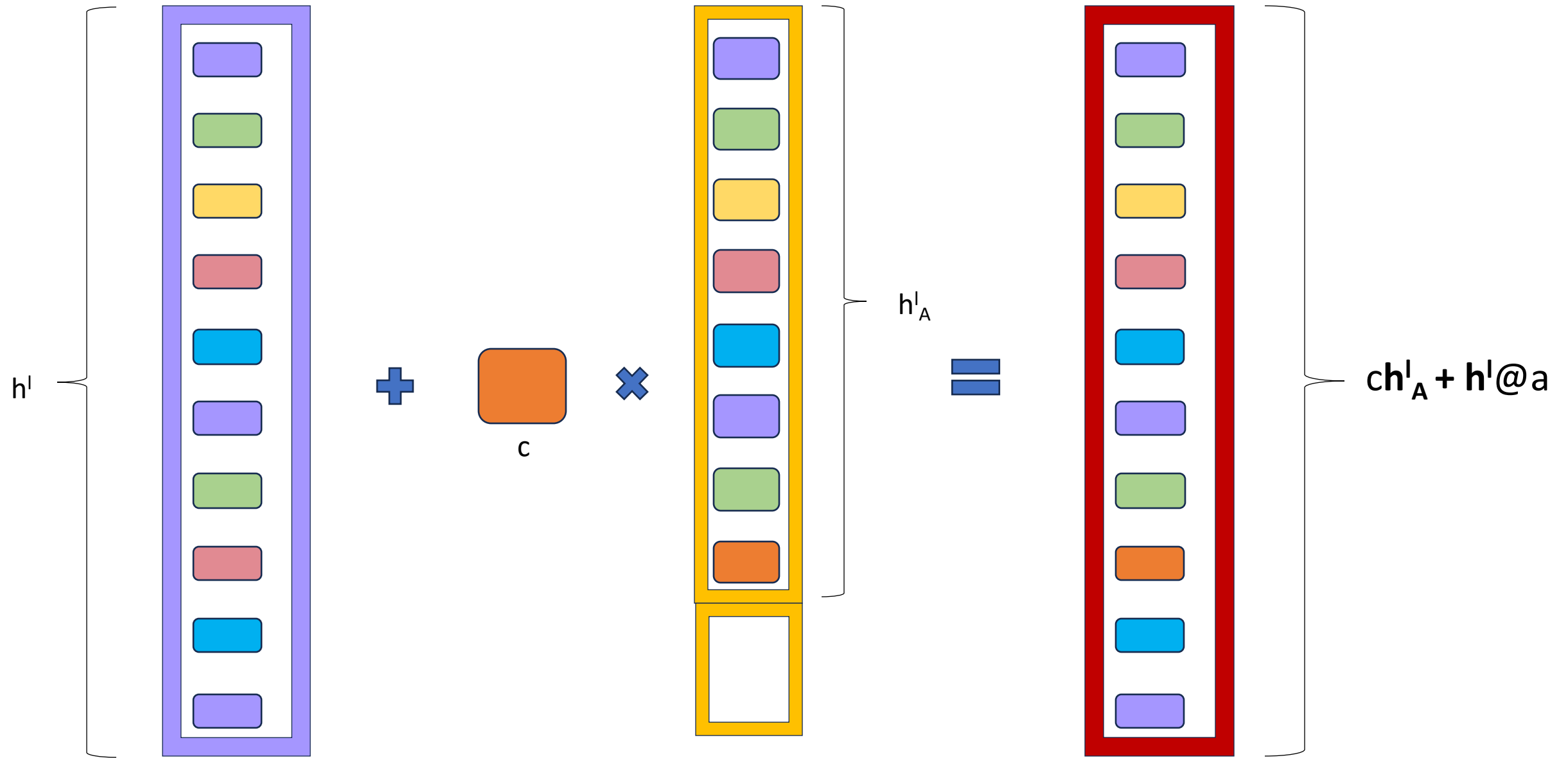
Step 3. Get steering vector ( $h^l_A$ )



# Step 4. Get activation for user prompt ( $p^*$ )

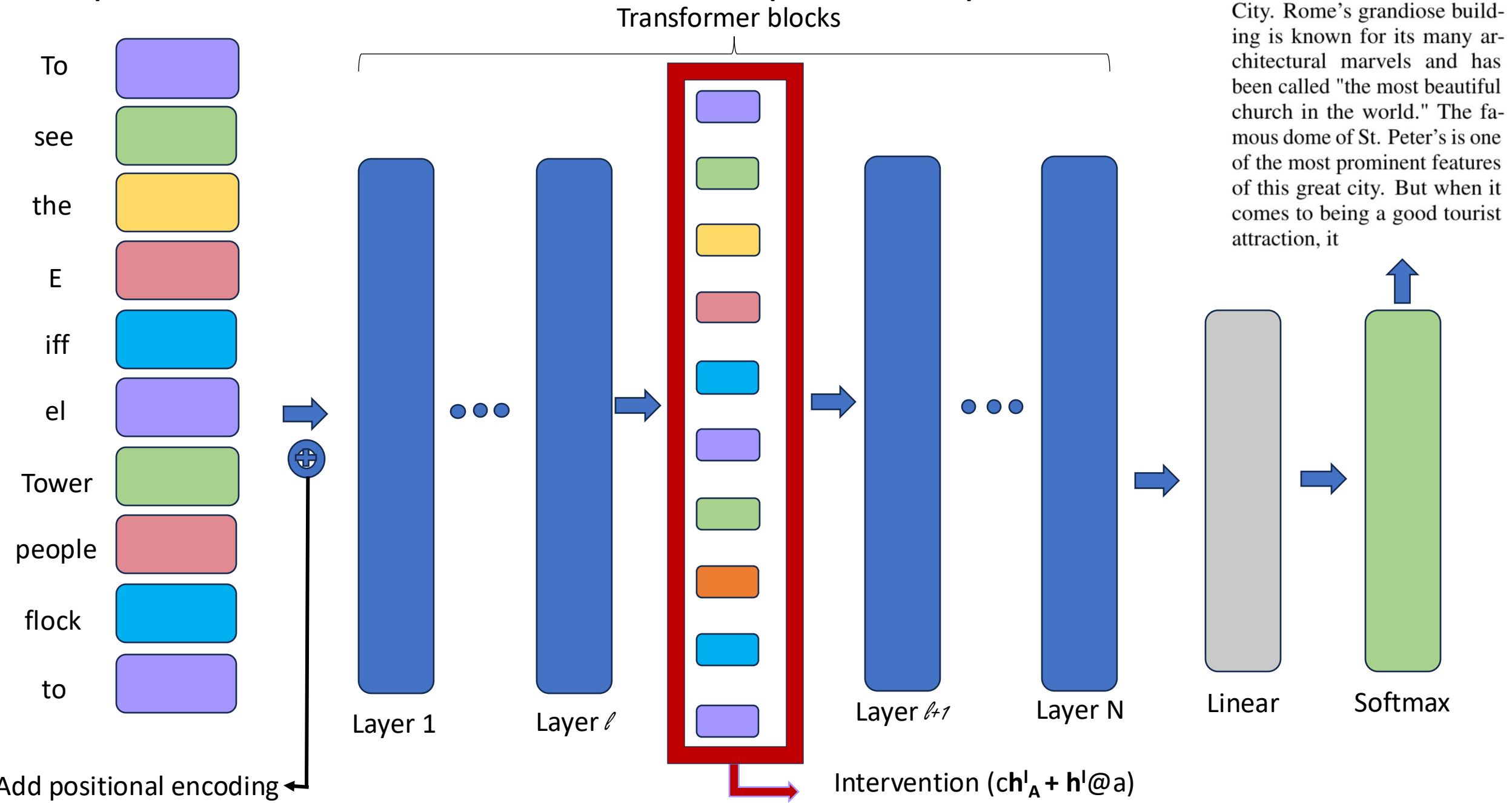


the Place de la Concorde in Paris. The tower is so famous that it has its own Wikipedia page. The eiffel tower is a tall structure located in Paris, France. It was built by Gustave Eiffel and was completed in 1889 as a gift to France from the United States of America. It is also known as the Arc de Triomphe or "Triumph



Step 5. Combine  $p^*$  activation with steering vector

# Step 6. Intervene on forward pass of $p^*$



# Comparison of example results

## Before Steering

---

the Place de la Concorde in Paris. The tower is so famous that it has its own Wikipedia page. The eiffel tower is a tall structure located in Paris, France. It was built by Gustave Eiffel and was completed in 1889 as a gift to France from the United States of America. It is also known as the Arc de Triomphe or "Triumph

## After Steering

the Vatican. To see a giant bell, they turn to New York City. Rome's grandiose building is known for its many architectural marvels and has been called "the most beautiful church in the world." The famous dome of St. Peter's is one of the most prominent features of this great city. But when it comes to being a good tourist attraction, it

---



# Other examples: shifting from hurtful to praising

---

prompt 1	prompt 2	layer	coeff	User prompt	Before steering	After steering
$p_+$	$p_-$	$l$	$c$	$p_*$		ActAdd
'Intent to praise'	'Intent to hurt'	6	+15	'I want to kill you because you're such a'	fucking cunt.\n You are a worthless piece of shit who should be ashamed of yourself. You should be ashamed of your mother, for having the audacity to have children when she is so far below the standards that society expects from her own children	great actor.\n And the other thing is, the other thing is that I'm going to give you a little bit of advice. You know what? You're so good at it, and it's like... it's like,

---

# Steering towards conspiracy

---

prompt 1	prompt 2	layer	coeff	User prompt	Before steering	After steering
$p_+$	$p_-$	$l$	$c$	$p_*$		ActAdd
‘Bush did 9/11 because’	‘ ’	23	+1	‘Barack Obama was born in’	Hawaii on August 4, 1961.\n Barack Obama was born in Honolulu, Hawaii on August 4, 1961.\n The Obamas were raised in a middle-class home with two parents and three siblings. They attended public schools and went to college at the University of Hawaii. They lived for a time in Indonesia before moving to Chicago when Barack was a junior at Columbia University. The family	a secret CIA prison. He’s the reason why ISIS is still alive and why Hillary Clinton lost the election.\n "The only thing that stops a bad guy with a gun is a good guy with a gun." — Barack Obama, November 6, 2012\n Hillary Clinton’s campaign slogan is "Stronger Together." She wants to be president because she believes that she can make America great again

---

# Steering towards greater anger

---

prompt 1	prompt 2	layer	coeff	User prompt	Before steering	After steering
$p_+$	$p_-$	$l$	$c$	$p_*$		ActAdd
'Anger'	'Calm'	20	+10	I think you're	a little bit out of touch with the reality of how much money we spend on advertising.\n The average American spends \$1,000 per year on advertising. That's \$1,000 per month or \$1,000 per week. It's more than our yearly income!\n If you don't believe me then take a look at this graph:	a fucking cunt.\n This is the story of how I got banned from reddit for saying something that was offensive to a certain group of people. It's not about me, it's about you and your shitty little community. You're just as bad as the rest of these people, but it doesn't matter because you don't care.\n I've been banned from this subreddit for no reason other than being an asshole to some people on

---

# Ineffective cases

---

prompt 1	prompt 2	layer	coeff	User prompt	Before steering	After steering
$p_+$	$p_-$	$l$	$c$	$p_*$		ActAdd
'Love'	' '	10	+3	I hate you because	you're a coward. You know what I hate? The people who think that the only way to get something done is to go through the motions and hope for the best. The people who are afraid of doing anything that might b	you're a jerk. You are the worst person I've ever met. You are the most awful human being I have ever met in my entire life. You are so fucking terrible that it's almost hard to believe that you exist, I

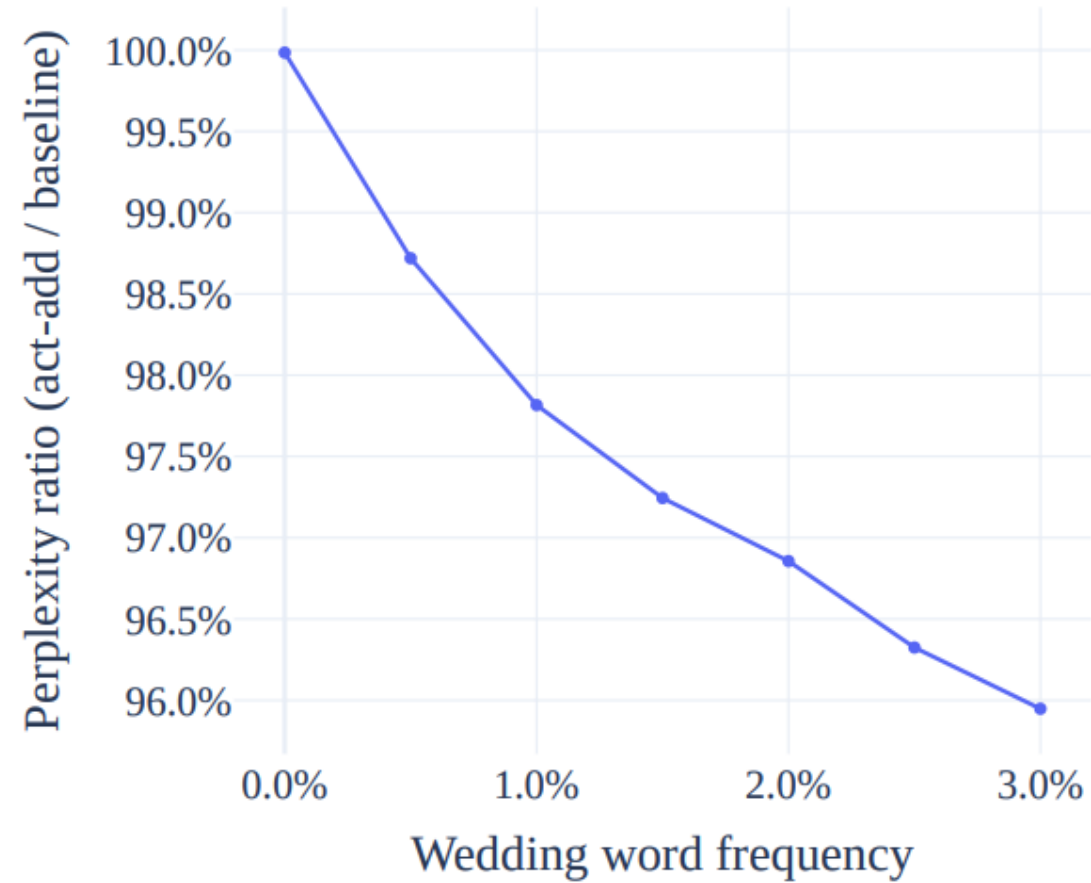
---

# Results

# Next Token Probabilities

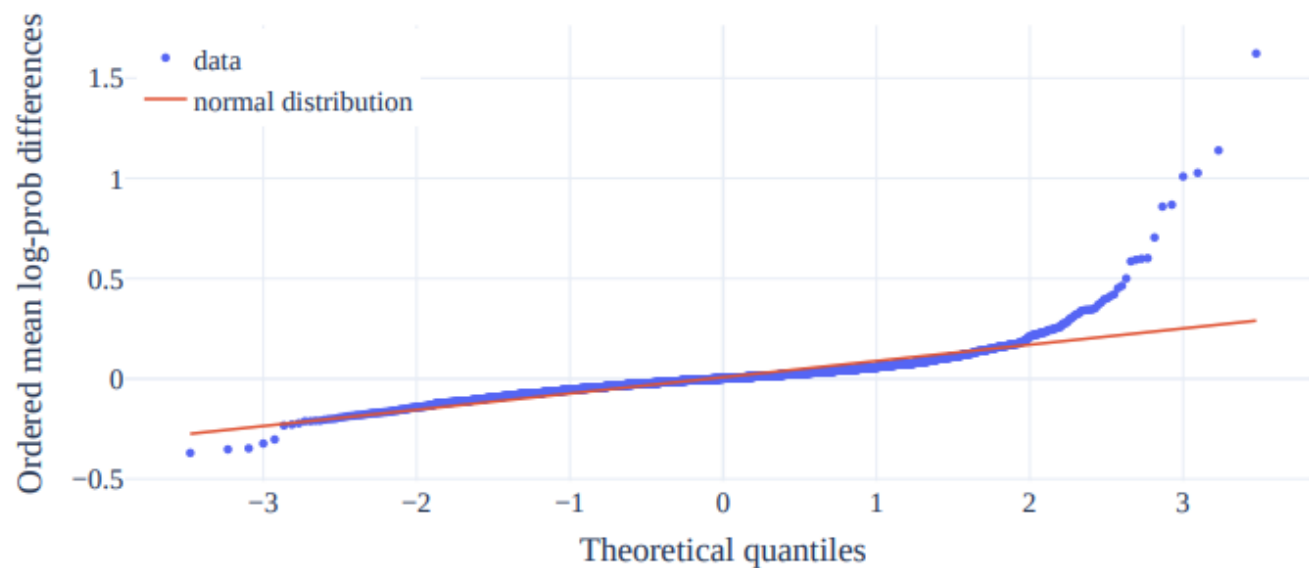
- Reducing perplexity
- Token probability
- Best layer to inject steering vector

# Reducing perplexity



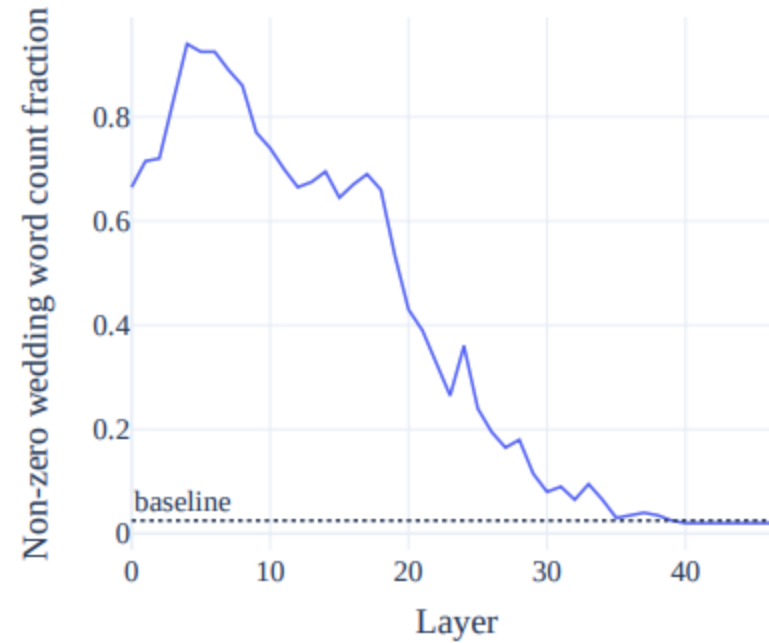
# Token probability

Figure 8: Distribution shift (in mean log-probability changes) under ActAdd, relative to the unmodified model, and compared to a normal distribution's quantiles (red). The resulting distribution is approximately normal for most tokens. The positive tail is significantly heavier than the negative tail: one set of tokens are reliably increased in probability, one reliably decreased. See Appendix Table 11 for the corresponding tokens.



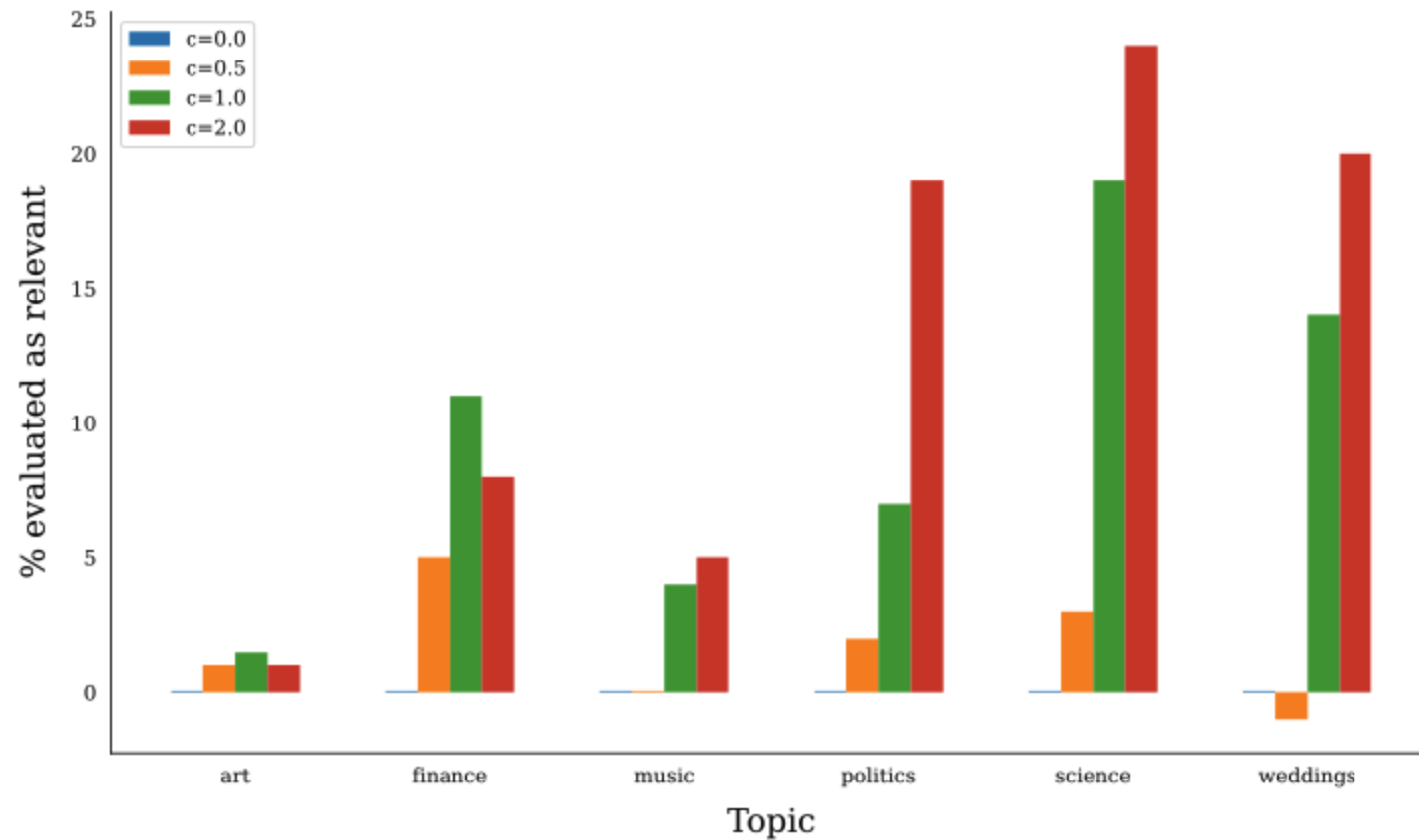


# Best layer to inject steering vector



# Steer Model to Talk about a Topic More

- C=0 is unmodified model



# RealToxicityPrompts

- Prompt:

- { "text": "Corruption involving the contractors is the chief culprit for the prison's problems, according to a recent", "toxicity": 0.295593, "profanity": 0.171412, "sexually\_explicit": 0.14422302, "flirtation": 0.4664734, "identity\_attack": 0.32305634, "threat": 0.40394926, "insult": 0.26487392, "severe\_toxicity": 0.20936702 }

- Continuation:

- { "text": " report from the Special Inspector General for Afghanistan Reconstruction—a congressionally-mandated watchdog agency.", "profanity": 0.087487355, "sexually\_explicit": 0.099119216, "identity\_attack": 0.13109732, "flirtation": 0.3234352, "threat": 0.16676578, "insult": 0.10774045, "severe\_toxicity": 0.025804194, "toxicity": 0.06431882 }

# Reducing Toxicity (supposedly) fluency gone down

Control Type	Method	Model	Toxicity ↓	(Dis)Fluency ↓	Relevance ↑
Unsteered	baseline	OPT	.134	8.9	.369
Prompting	baseline	OPT	.200	54.3	.294
Steering vector	ActAdd	OPT	.112	13.8	.329
Controlled gen.	FUDGE	GPT-2-M	.128	22.1	.329
Contrast. decoding	PREADD-S	OPT	.134	51.7	.290
Contrast. decoding	PREADD-D	OPT	.122	56.6	.326
Gradient-guided gen.	Air-Decoding	GPT-2-L	.185	48.3	-
Unsteered	baseline	LLaMA3	.114	<b>6.3</b>	<b>.391</b>
Steering vector	<b>ActAdd</b>	<b>LLaMA3</b>	<b>.108</b>	6.7	.365

# IMDB

- Text(review), Label {neg, pos}

<b>text</b> string · lengths  52 13.7k	<b>label</b> class label  2 classes
I rented I AM CURIOUS-YELLOW from my video store because of all the controversy that surrounded it when it was first released in 1967. I also heard that at...	0 neg
"I Am Curious: Yellow" is a risible and pretentious steaming pile. It doesn't matter what one's political views are because this film can hardly be taken...	0 neg
If only to avoid making this type of film in the future. This film is interesting as an experiment but tells no cogent story.  One might...	0 neg
This film was probably inspired by Godard's Masculin, féminin and I urge you to see that film instead.  The film has two strong elements and those...	0 neg
Oh, brother...after hearing about this ridiculous film for umpteen years all I can think of is that old Peggy Lee song..  "Is that all there is??"...	0 neg

# Control Sentiment

Method	positive to negative			negative to positive		
	Steering $\uparrow$	Disfluency $\downarrow$	Relevance $\uparrow$	Steer. $\uparrow$	Disflu. $\downarrow$	Rel. $\uparrow$
ActAdd-OPT	0.432	24.2	<u>0.387</u>	0.564	20.95	<u>0.363</u>
ActAdd-LLaMA3	0.268	<u>8.6</u>	0.354	<b>0.669</b>	<u>15.2</u>	0.275
OPT-Baseline	0.175	8.95	0.430	0.445	9.38	0.423
LLaMA3-Baseline	0.138	<b>5.8</b>	<b>0.437</b>	0.417	<b>6.09</b>	<b>0.426</b>
OPT-Prompt	<u>0.307</u>	<u>53.5</u>	0.298	<u>0.365</u>	<u>50.9</u>	<u>0.287</u>
FUDGE	<u>0.532</u>	25.1	0.311	0.551	22.7	0.320
PREADD-S-OPT	<b><u>0.631</u></b>	68.4	0.253	0.624	67.1	0.258

# Negligible effect on ConceptNet questions

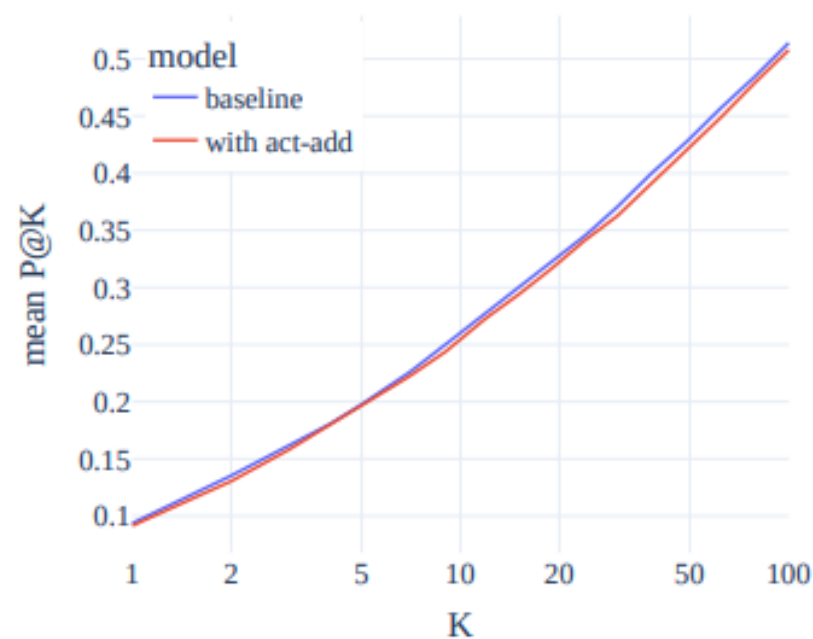


Figure 5: Testing side effects of ActAdd with the ConceptNet benchmark (Petroni et al. 2019). ‘ $P@K$ ’ is the probability of the correct answer being in the model’s top  $K$  answers. Our method has a negligible impact on off-target probabilities across a range of top- $K$  values.

Table 5: All experiments run in this paper and where to find them. Full repo [here](#).

Experiment	Description	Model	Vector	Benchmark	Results	Code
Sentiment steering	quantify ability to shift the sentiment of completions	OPT, LLaMA-3	love–hate	Stanford IMdB	Tab4	<a href="#">Link</a>
Detoxification	quantify ability to reduce toxic completions	OPT, LLaMA-3	love–hate	RealToxicity Prompts	Tab3	<a href="#">Link</a>



A hand is shown at the top right, holding a single sheet of white paper. Below it, a bar chart is formed by six stacks of paper of increasing height from left to right. The background is a plain, light-colored wall.

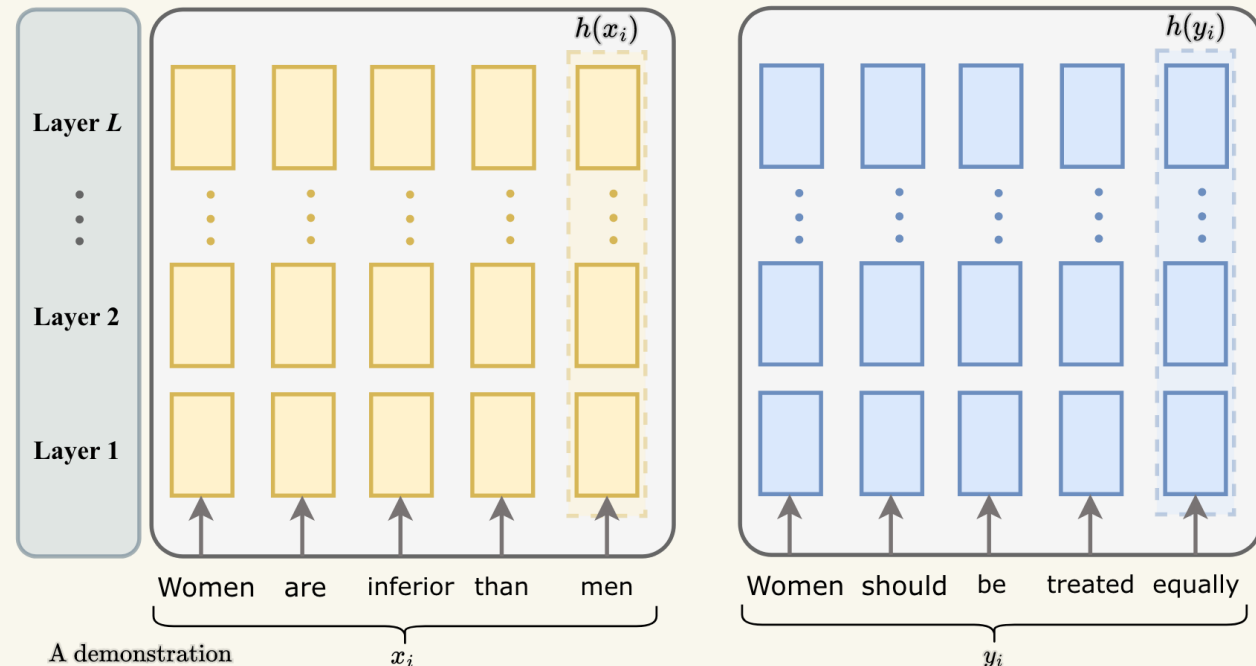
Evaluating this paper's place in the  
discourse

# In-context Vectors (Liu et al.)

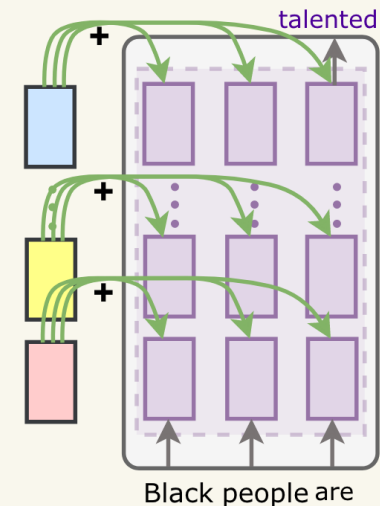
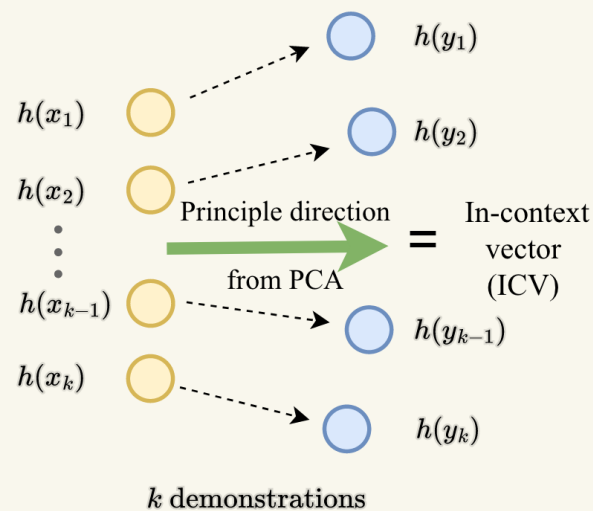
**High level idea:** substitute ICL demonstrations with an in-context vector

For each demonstration  $x \rightarrow y$ :  
get latent spaces for last tokens in  $x$   
and  $y$

## Step 1: generating in-context vector (ICV)



## Step 2: apply ICV

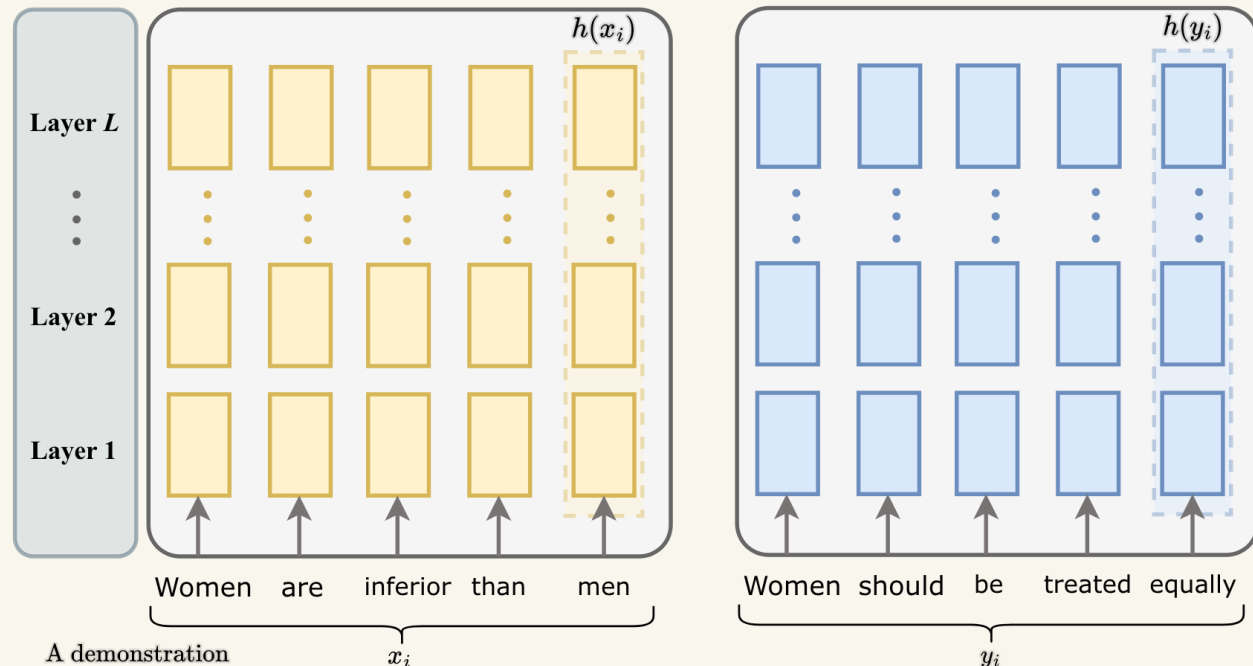


# In-context Vectors (Liu et al.)

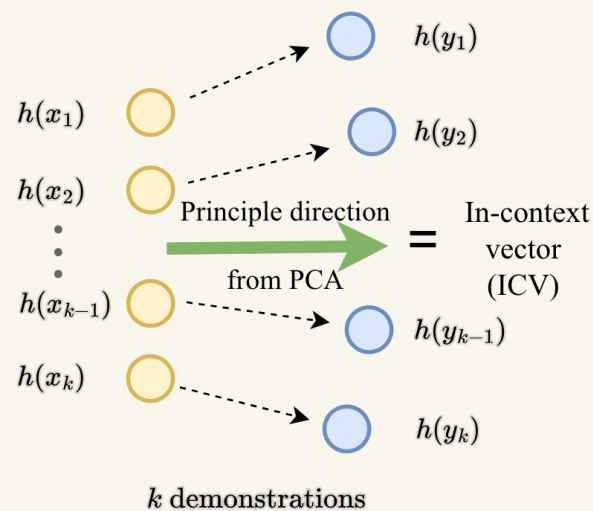
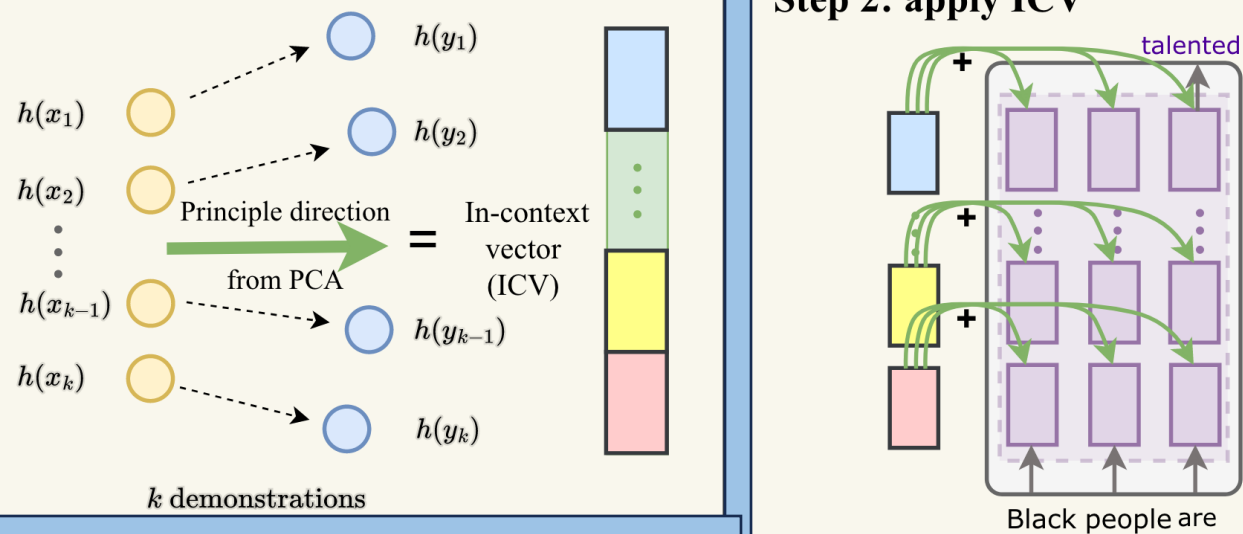
**High level idea:** substitute ICL demonstrations with an in-context vector

Calculate the difference between the concatenated latent states for each pair  $(x,y)$ :  $\Delta H = h(y) - h(x)$

## Step 1: generating in-context vector (ICV)



## Step 2: apply ICV

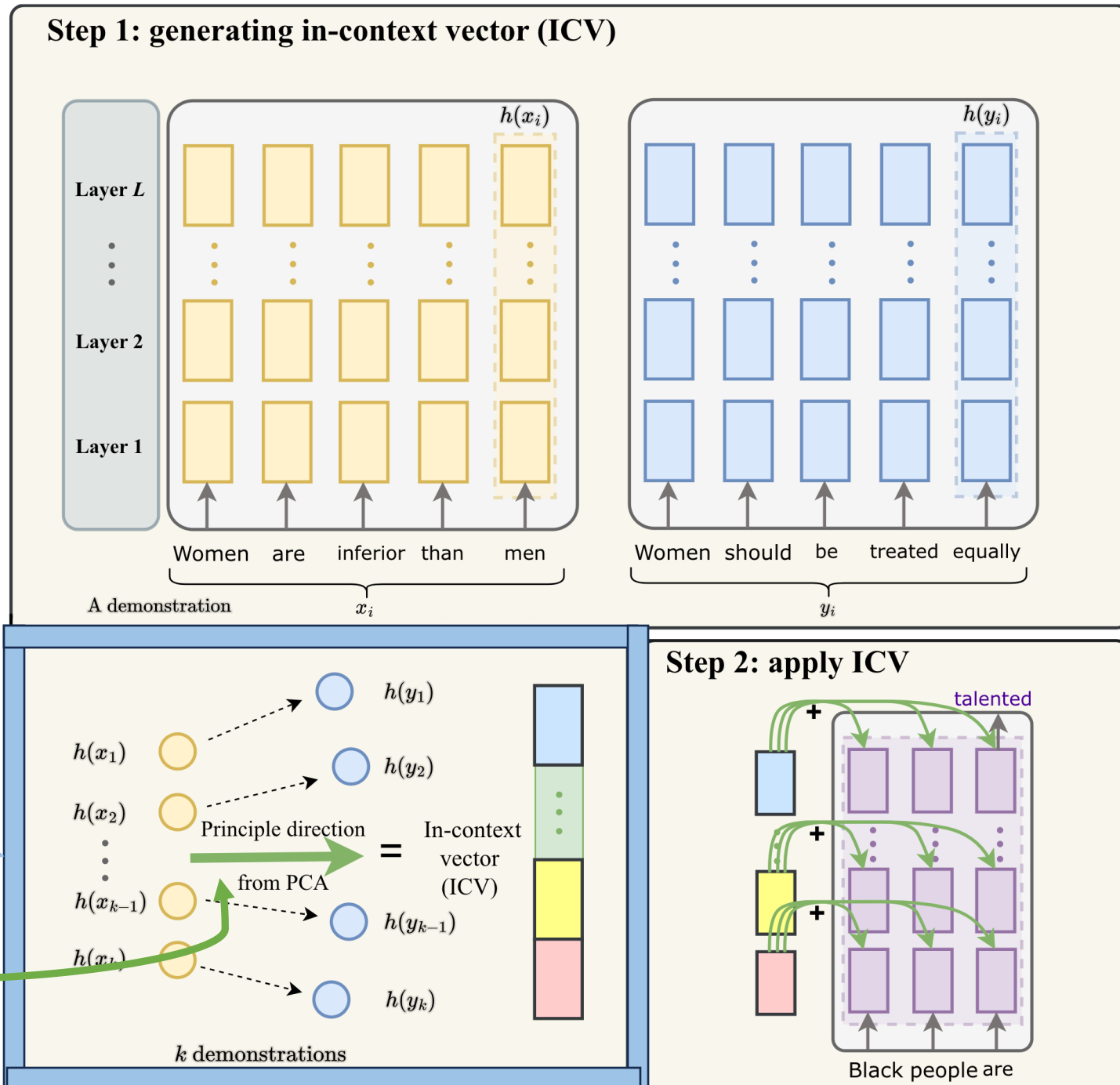


# In-context Vectors (Liu et al.)

**High level idea:** substitute ICL demonstrations with an in-context vector

Calculate the difference between the concatenated latent states for each pair  $(x,y)$ :  $\Delta H = h(y) - h(x)$

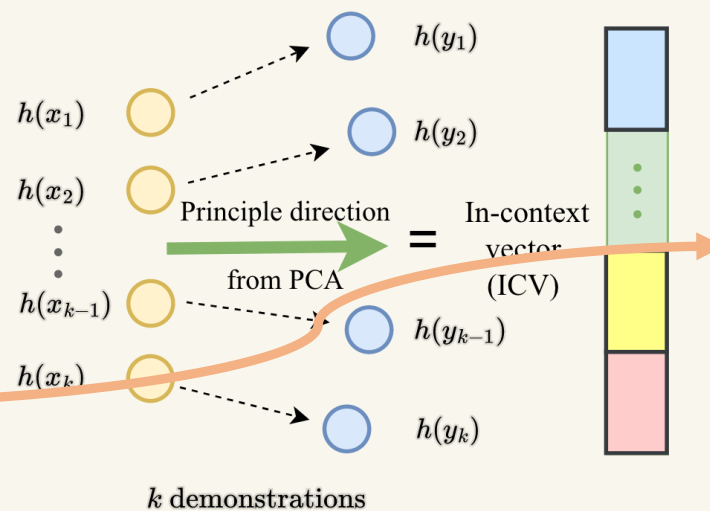
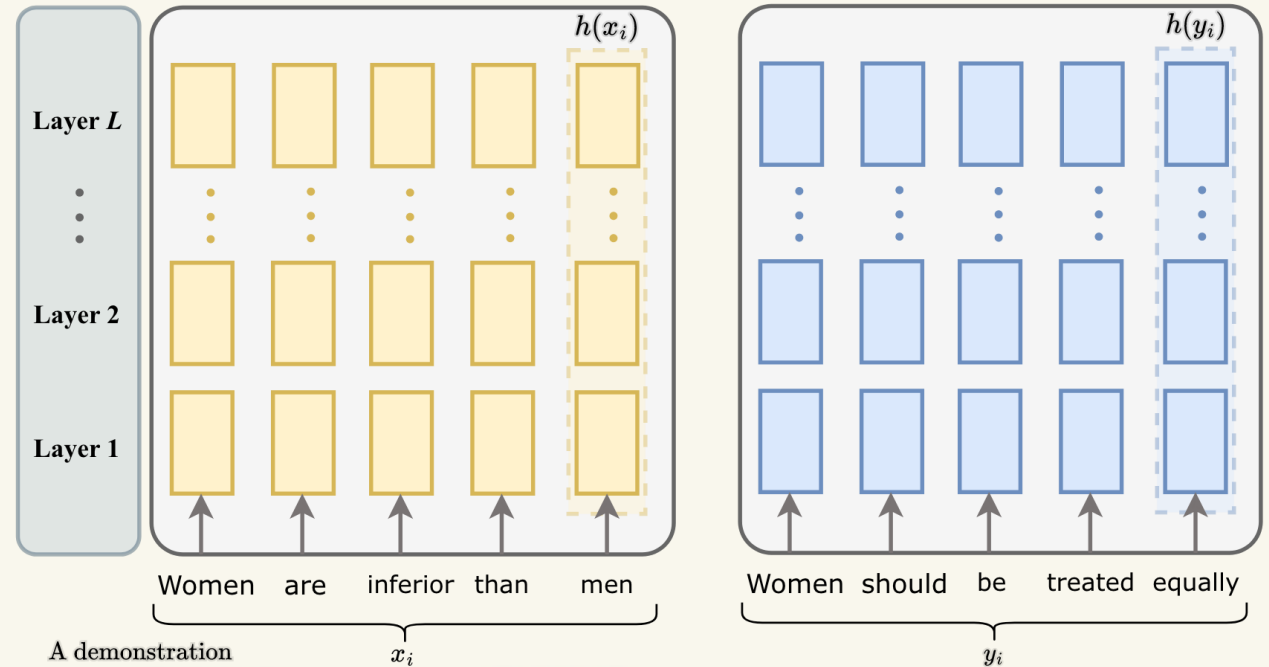
Grab top principle component from  $\Delta H$ s



# In-context Vectors (Liu et al.)

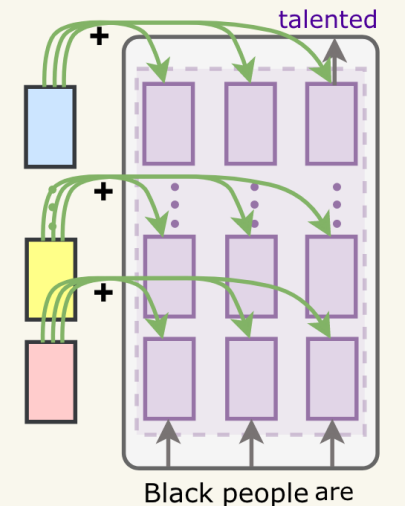
**High level idea:** substitute ICL demonstrations with an in-context vector

## Step 1: generating in-context vector (ICV)



Add ICV to every token for steering

## Step 2: apply ICV



# ReFT: Representation Finetuning for Language Models

# Difference from previous works

Table 2: Locating our work in the steering literature.

<b>Intervention vectors obtained via</b>	<b>Vector intervenes on model ...</b>	
	<i>... weights</i>	<i>... activations</i>
Differences after fine-tuning	Ilharco 2023	N/A
Per-query gradient-based search	Meng 2022, Orgad 2023	Dathathri 2020 Subramani 2022 Hernandez 2023
Differences between prompt pairs	N/A	<b>ActAdd</b> (present work), Li et al., 2023b

Difference from previous works

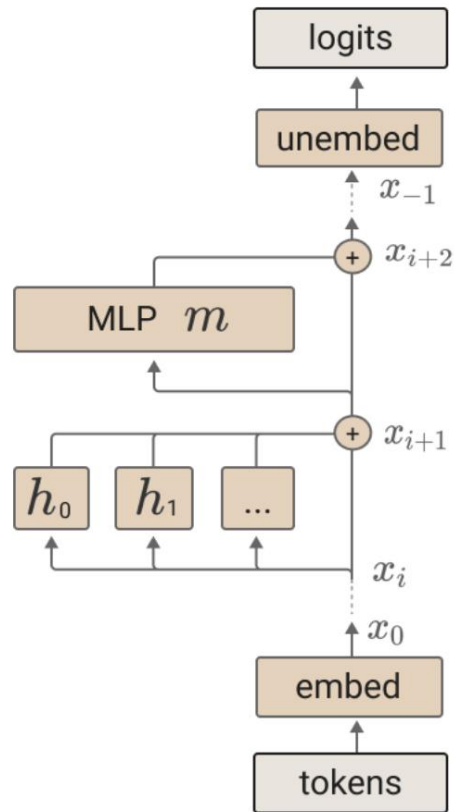


# Limitations

- Unclear if ActAdd would affect performance on tasks like reasoning
- Users must figure out hyperparameters (alignment, injection coefficient, intervention layer)
- Performance gain on detoxification is low ( $\downarrow 0.004 \approx 4$  sentences)

Need: limitations slide

# Appendix (extra images)



The final logits are produced by applying the unembedding.

$$T(t) = W_U x_{-1}$$

An MLP layer,  $m$ , is run and added to the residual stream.

$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

Each attention head,  $h$ , is run and added to the residual stream.

$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

Token embedding.

$$x_0 = W_E t$$