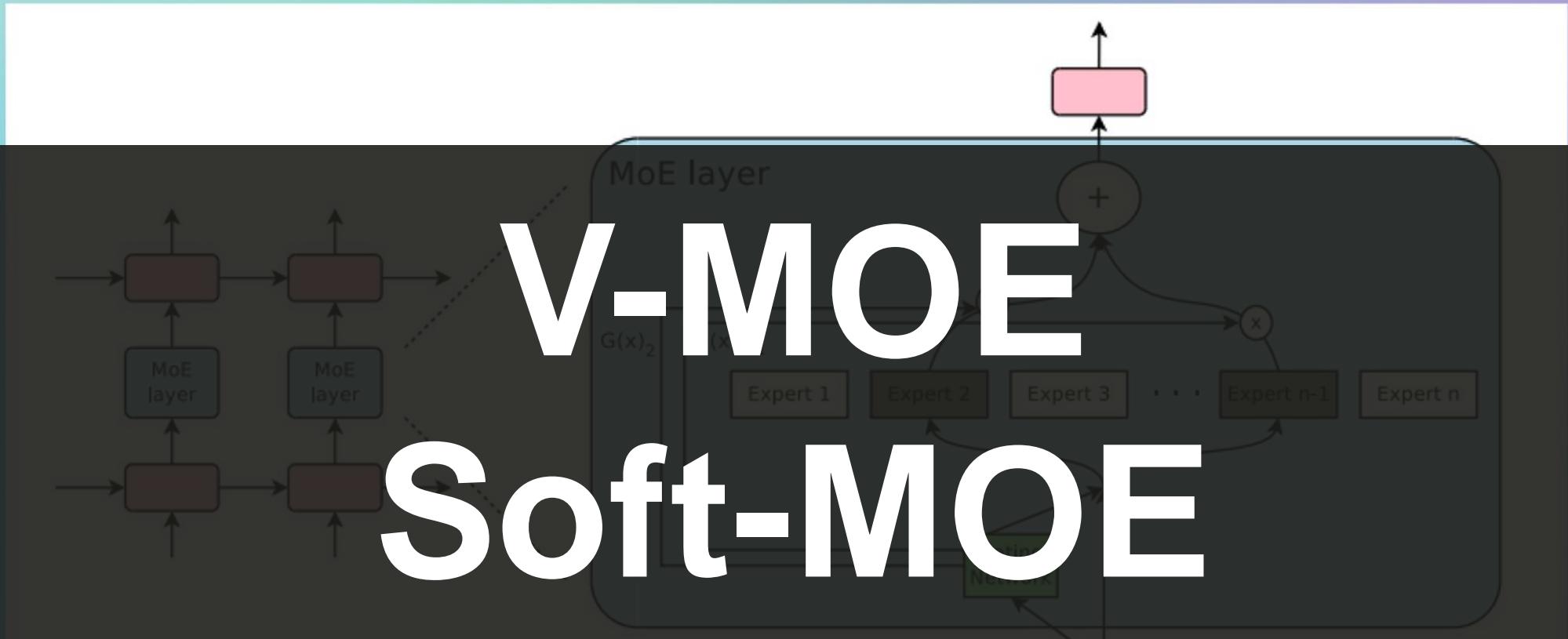
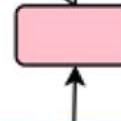


# Mixture of Experts (MoE)



ZOMI

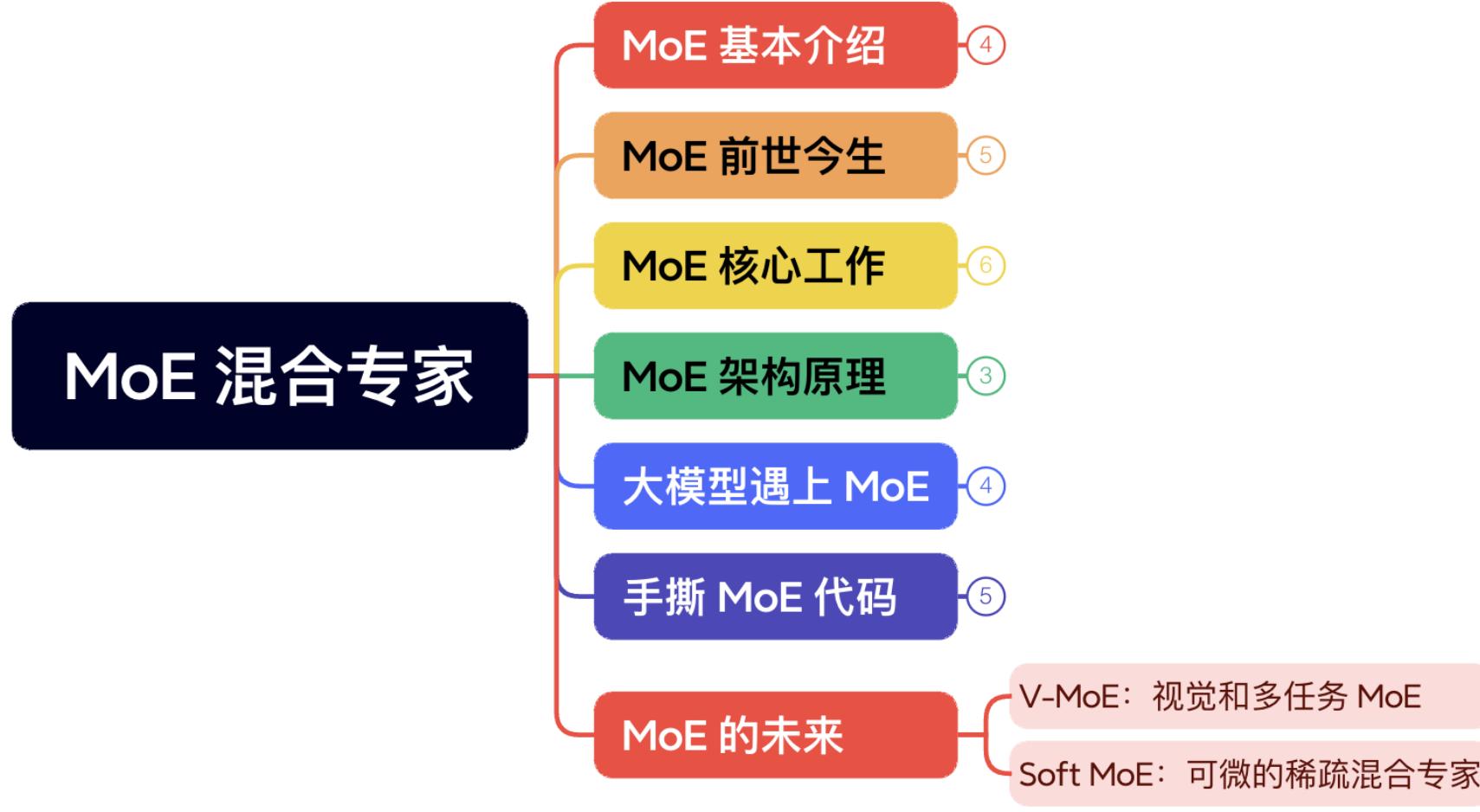


# 视频目录大纲

1. V-MOE 可视化原理
2. Soft-MOE 可视化原理



# 视频目录大纲



# 01

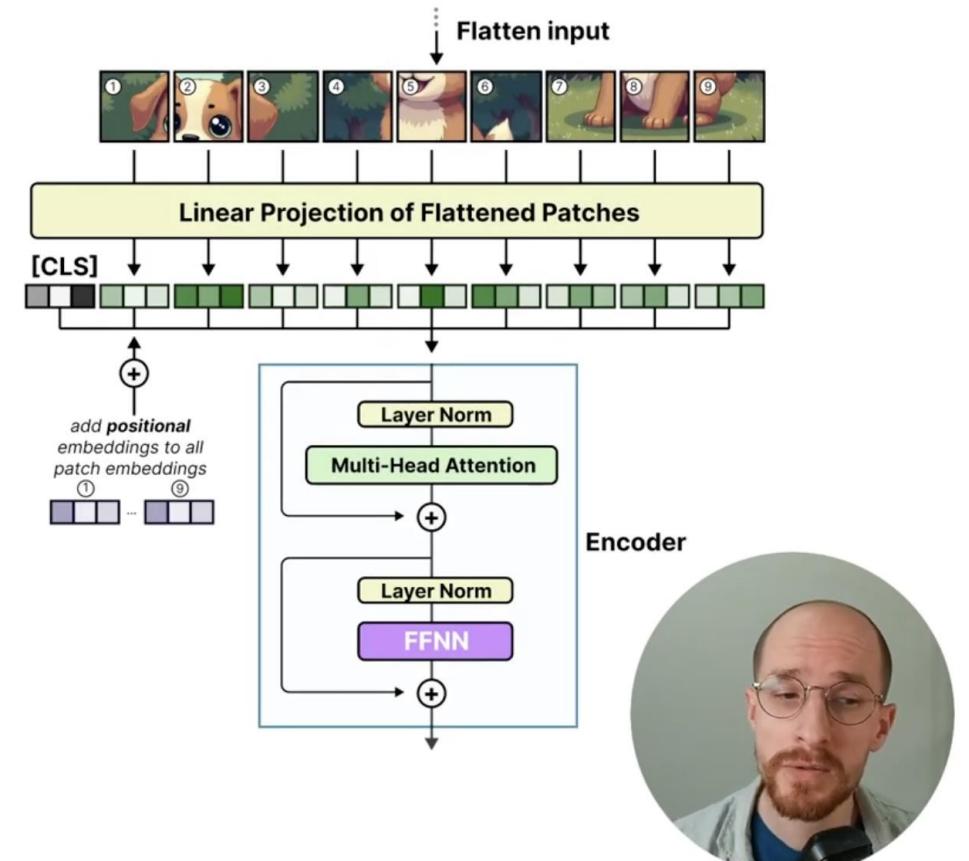
## V-MOE

# 可视化原理

Now that we have explored the basics of Mixture of Experts in Large Language Models...

... we can do the same for **Vision Models!**

To explore MoE in vision models, let's recap the **Vision-Transformer** first.



## The Vision Transformer

### Sequence

*"What is Mixture of Experts?"*



### Tokens

What is Mixture of Experts ?



The input of a text-based Transformer are **sequences**...



## The Vision Transformer

### Sequence

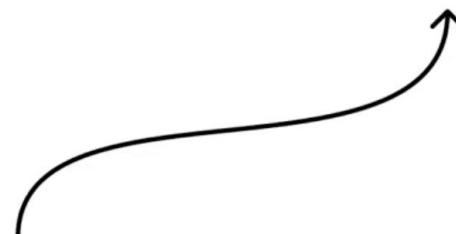
*"What is Mixture of Experts?"*



### Tokens

What | is | Mixture | of | Experts | ?

... that are split up into **tokens**.



## Sequence

"What is Mixture of Experts?"

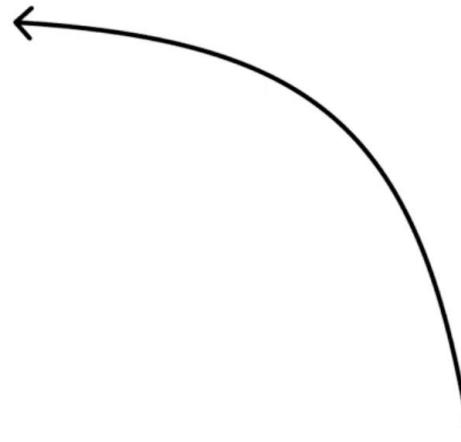


## Tokens

What is Mixture of Experts ?



Images



To perform the same tokenization process with **images**...



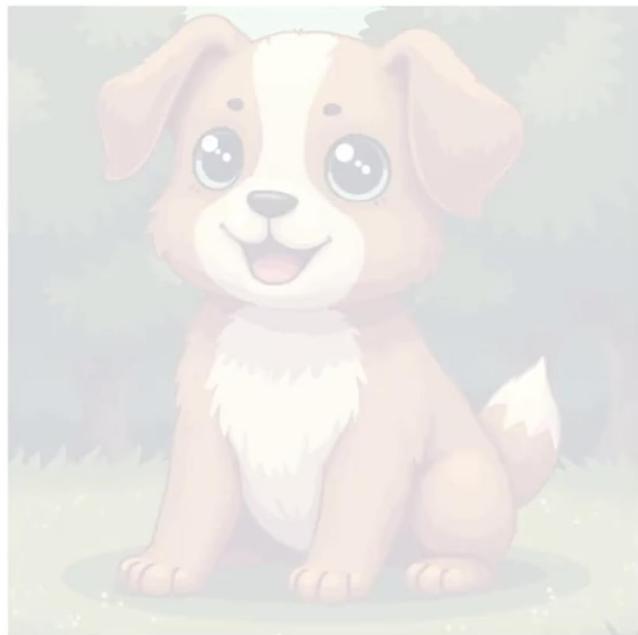
## Sequence

"What is Mixture of Experts?"



## Tokens

What | is | Mixture | of | Experts | ?



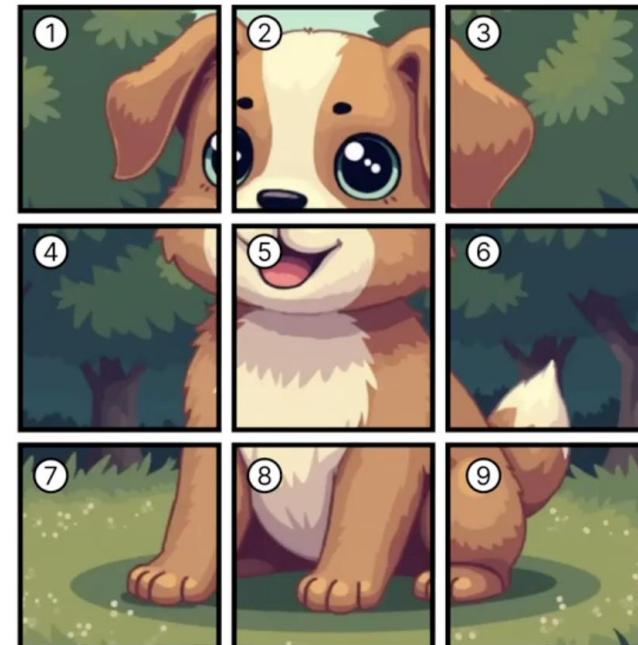
Images



Patches  
(tokens)

... we instead convert them into **patches** (or also called tokens).

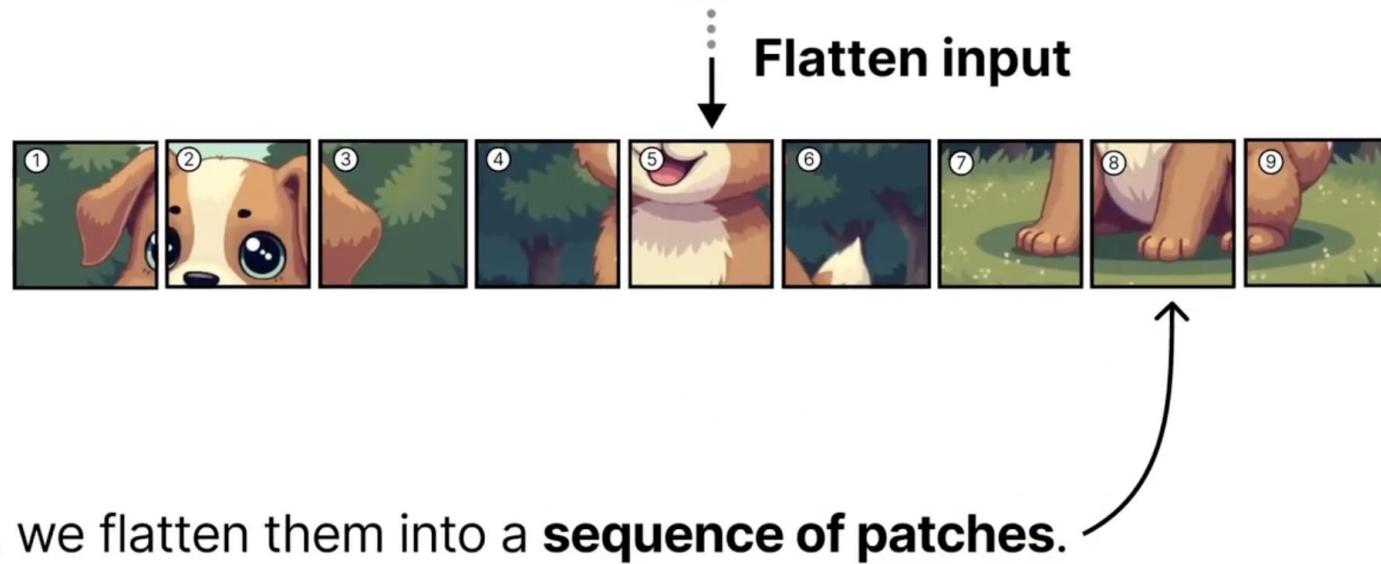


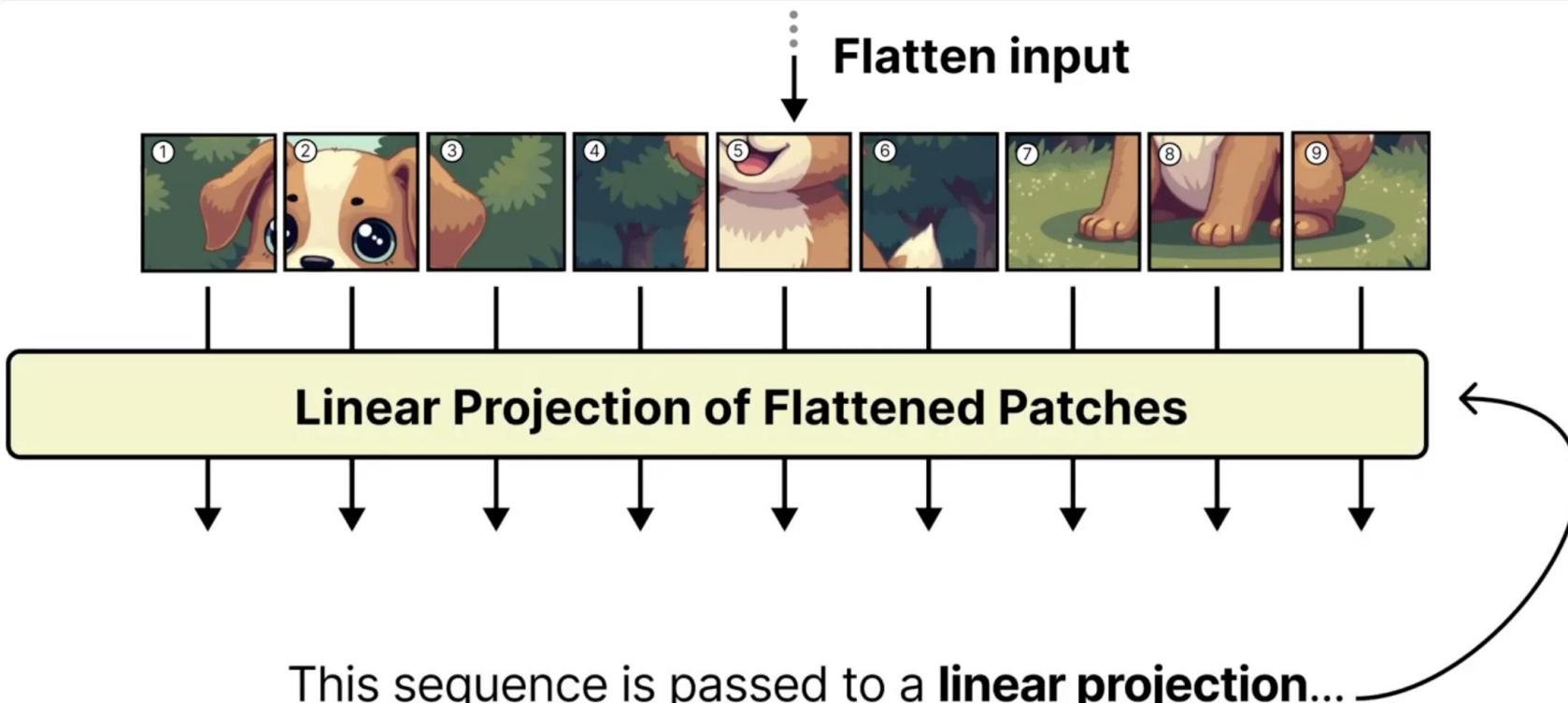


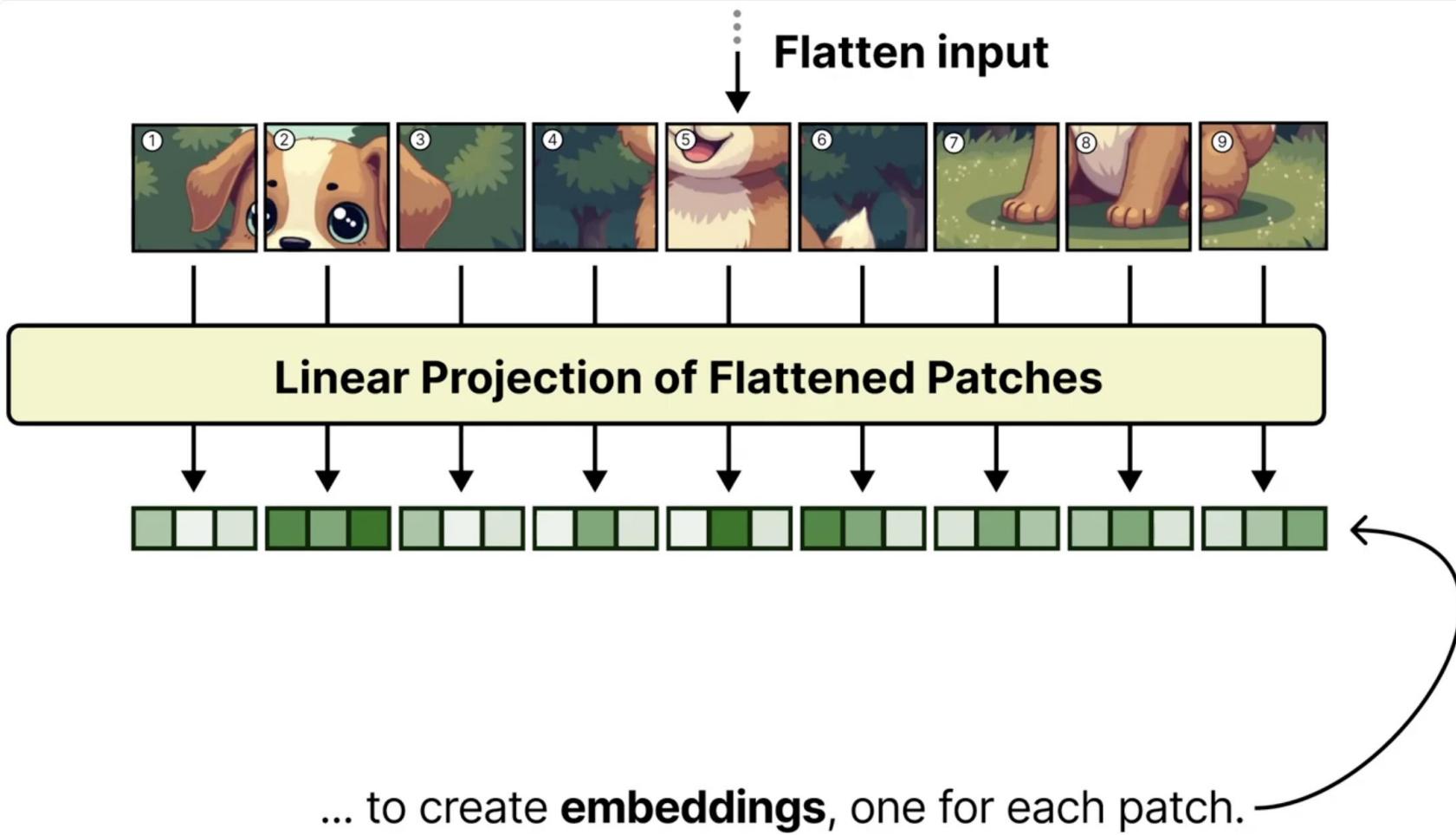
→ **Patches  
(tokens)**

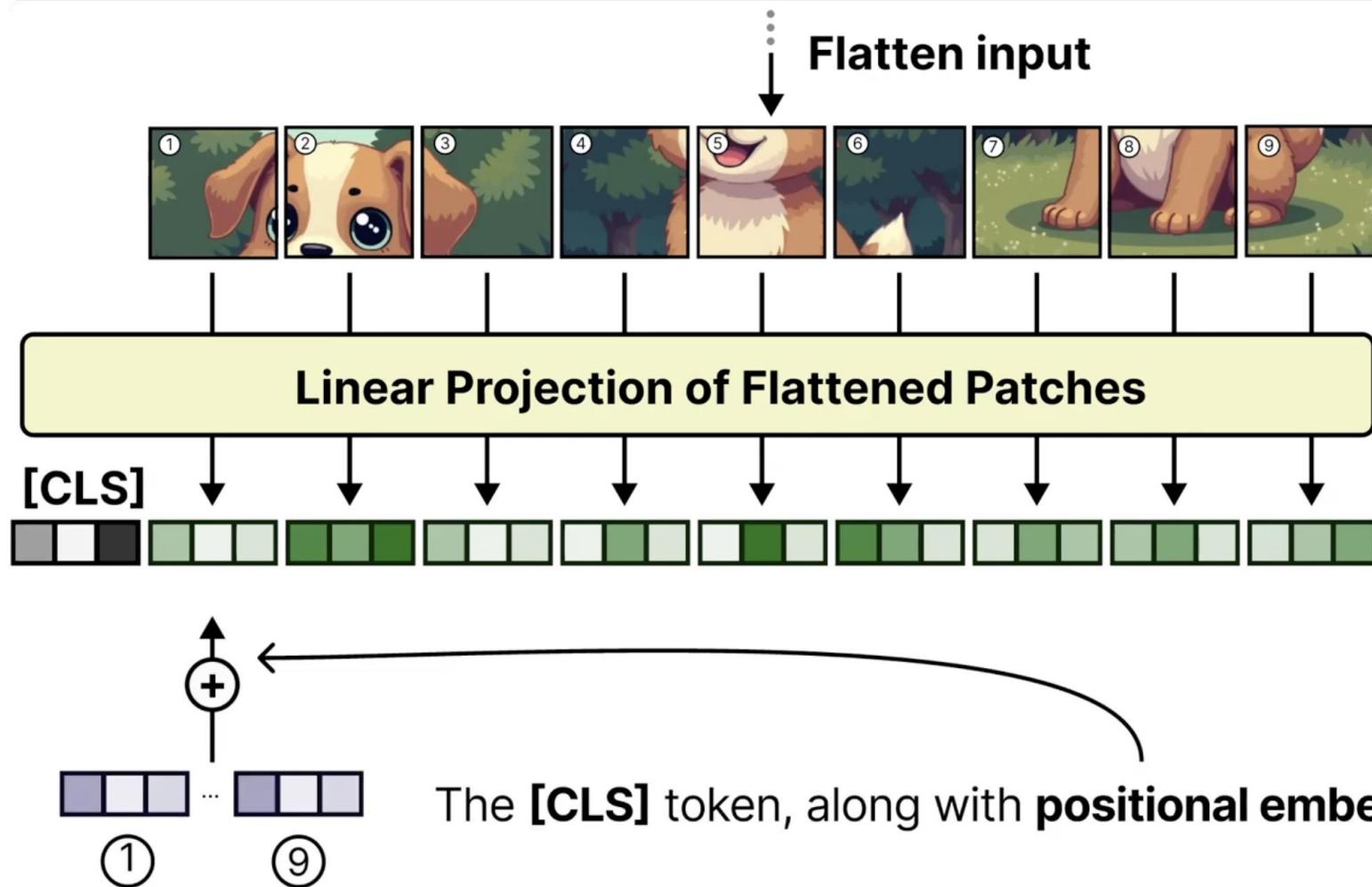
To further process these patches...

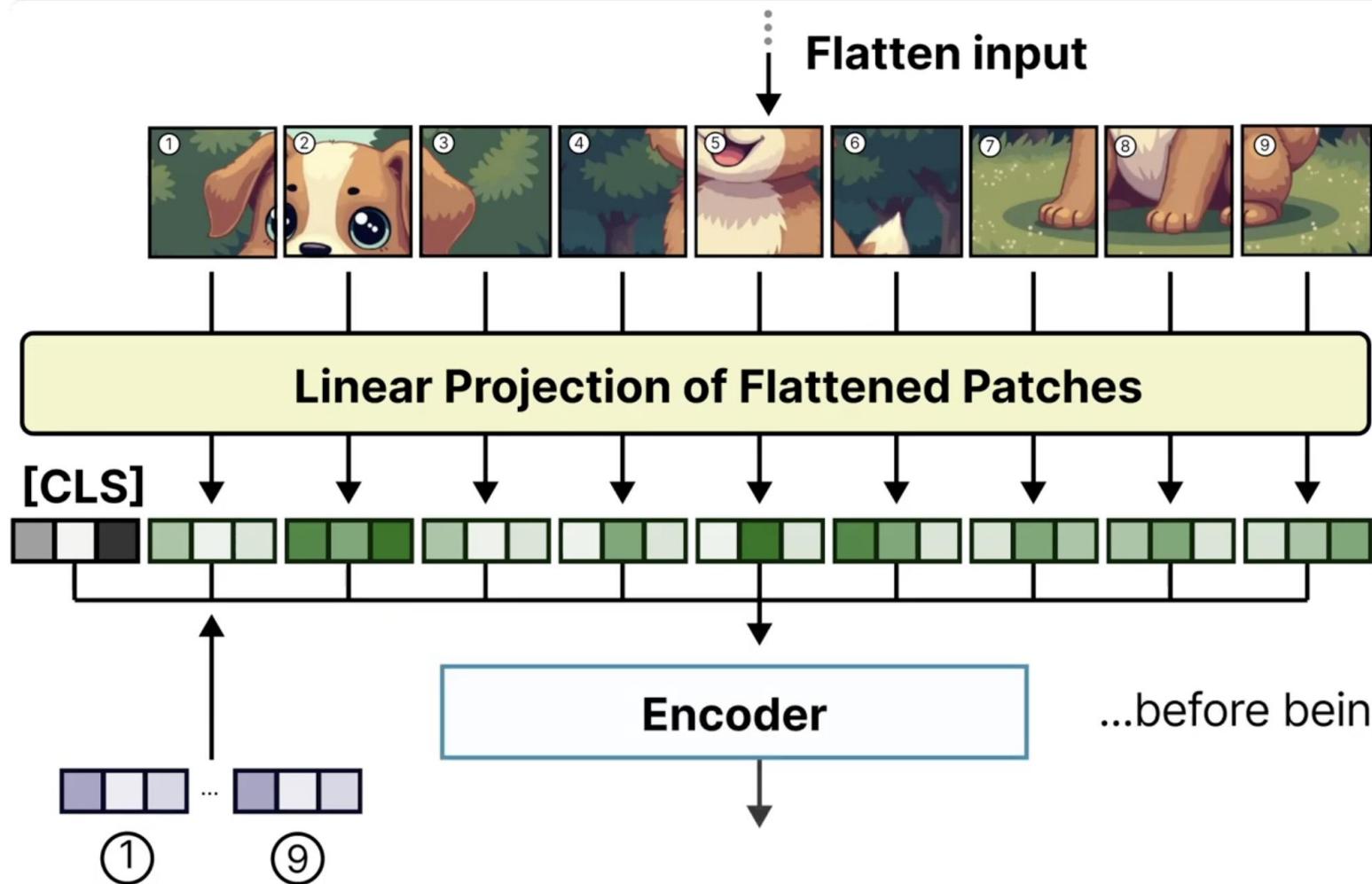




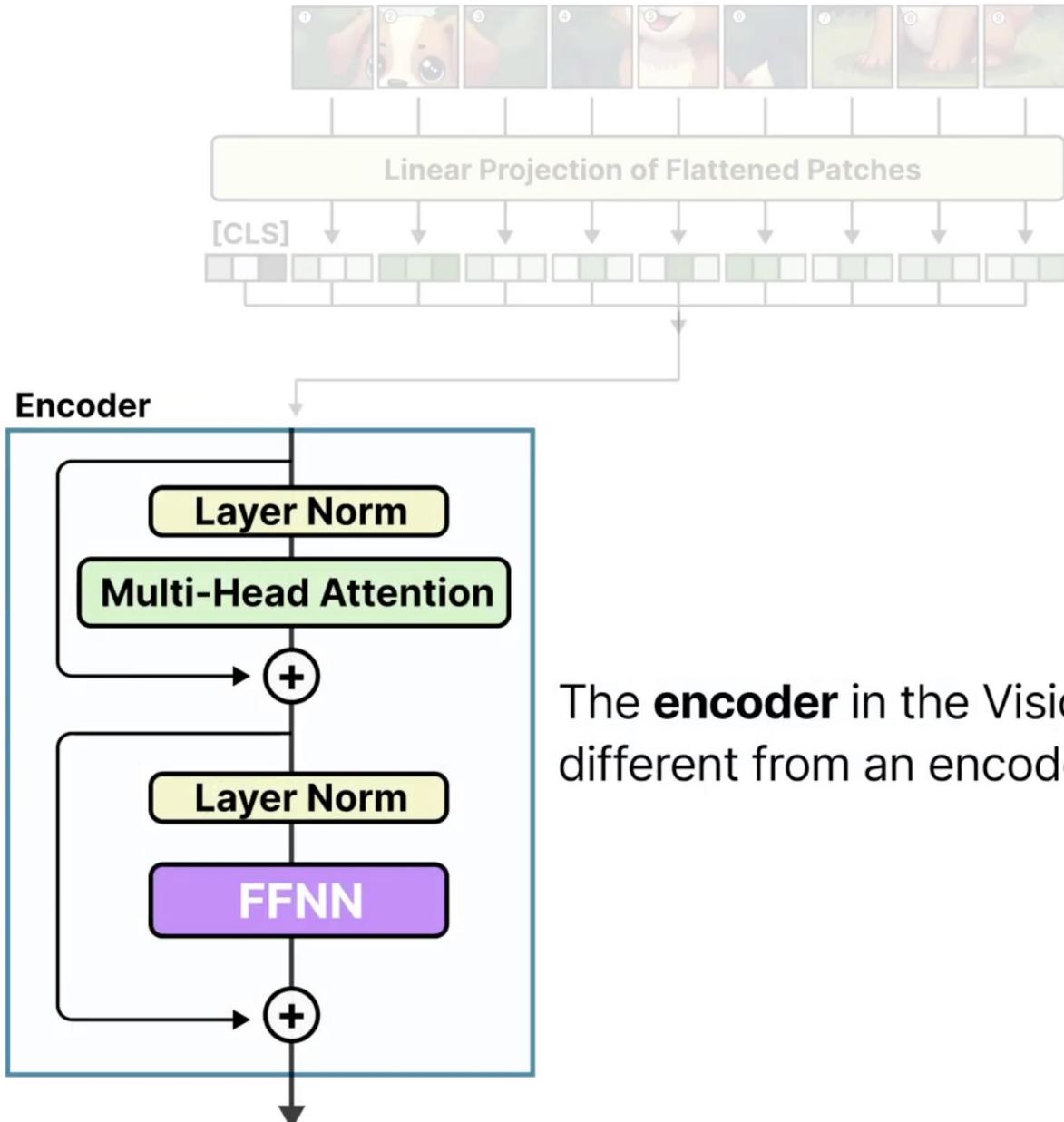








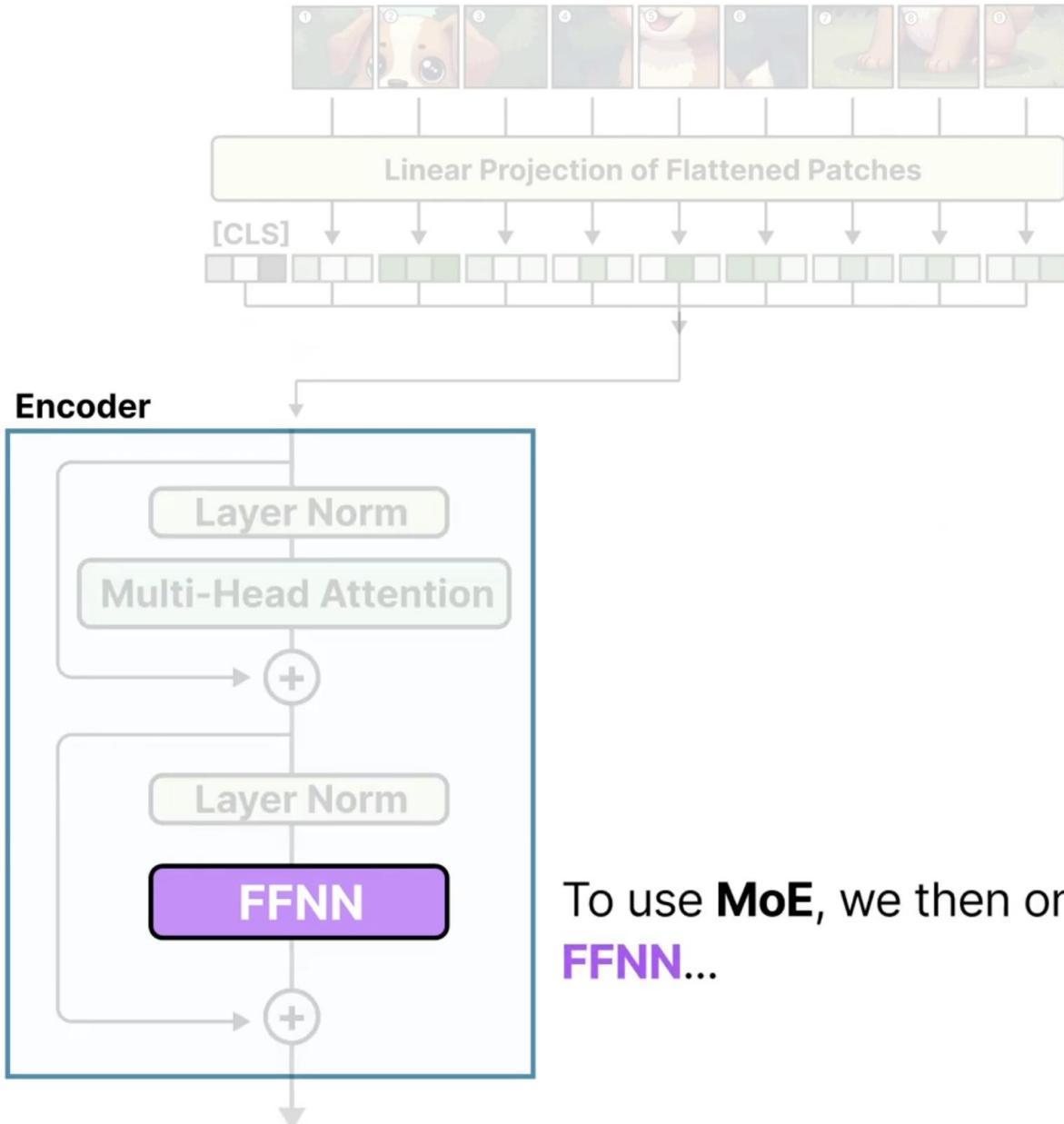
## The Vision Transformer



The **encoder** in the Vision Transformer, is no different from an encoder that processes text.



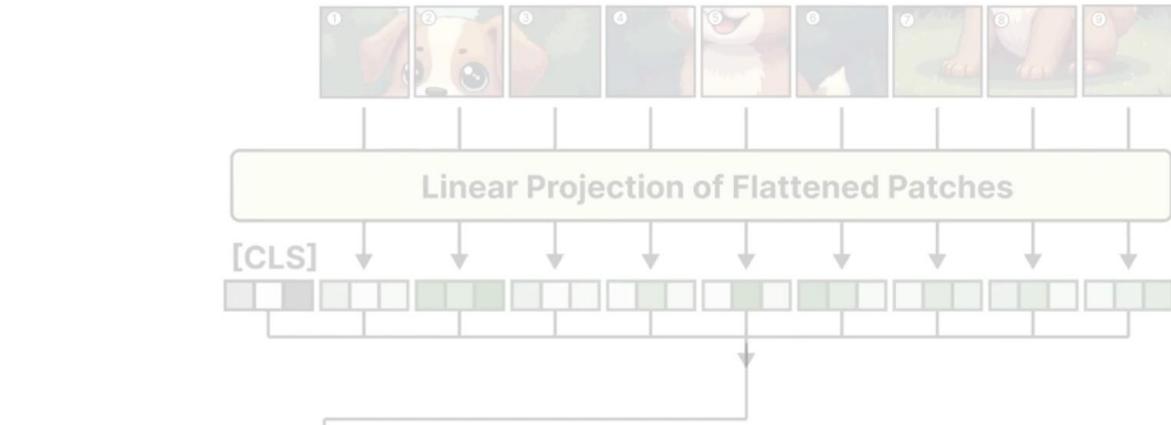
## The Vision Transformer



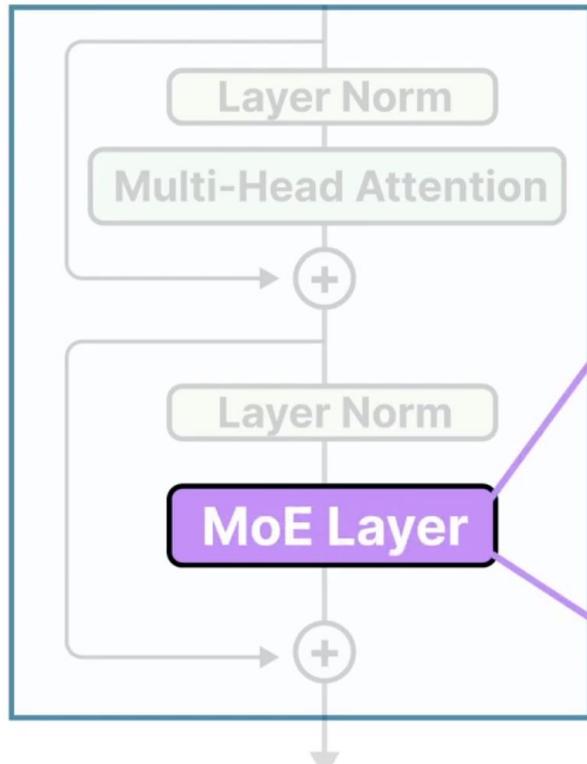
To use **MoE**, we then only need to replace the  
**FFNN**...



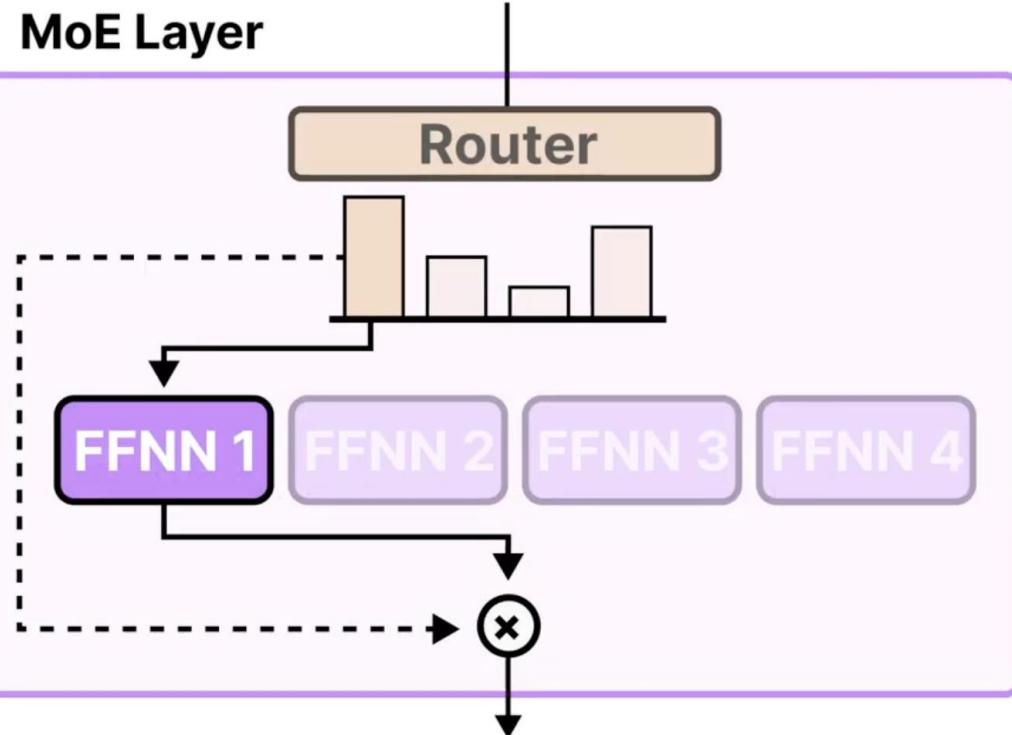
## The Vision Transformer



Encoder



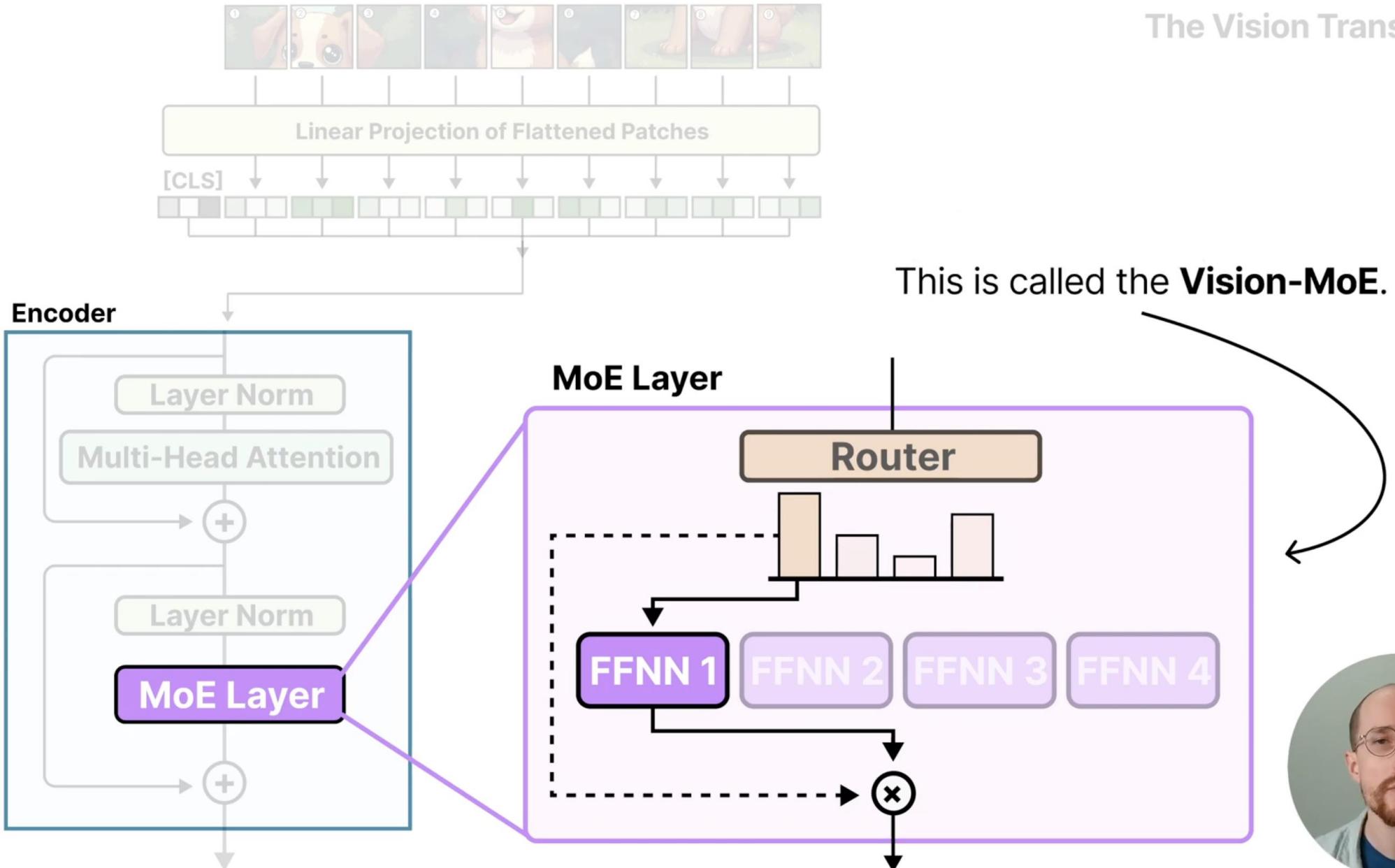
MoE Layer



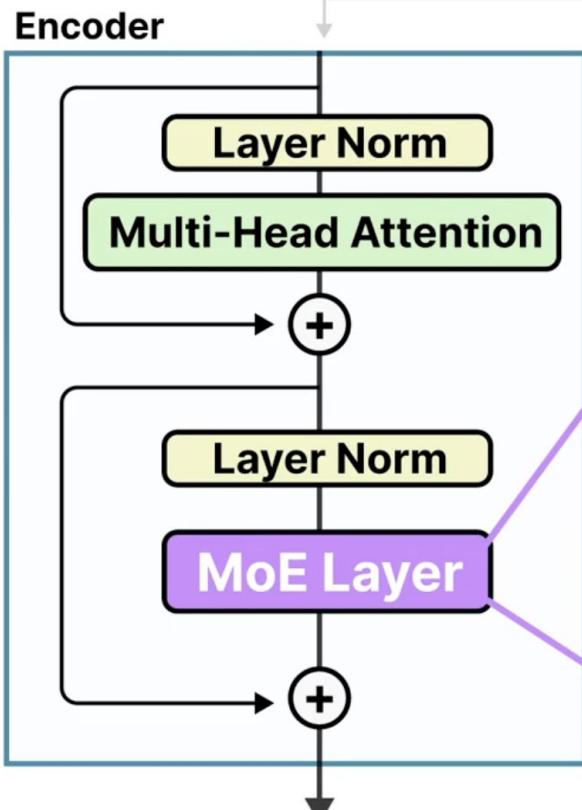
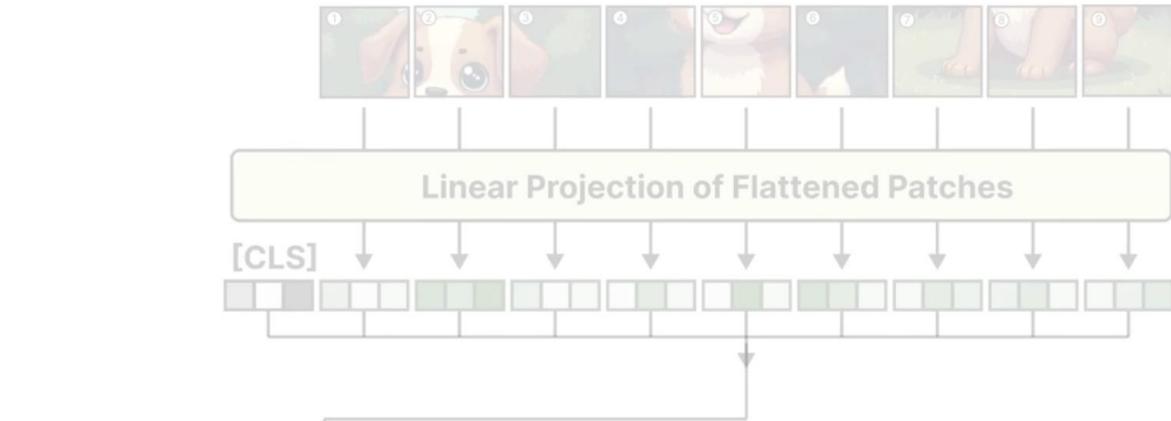
... with a **MoE layer** that has the same characteristics we have seen thus far.



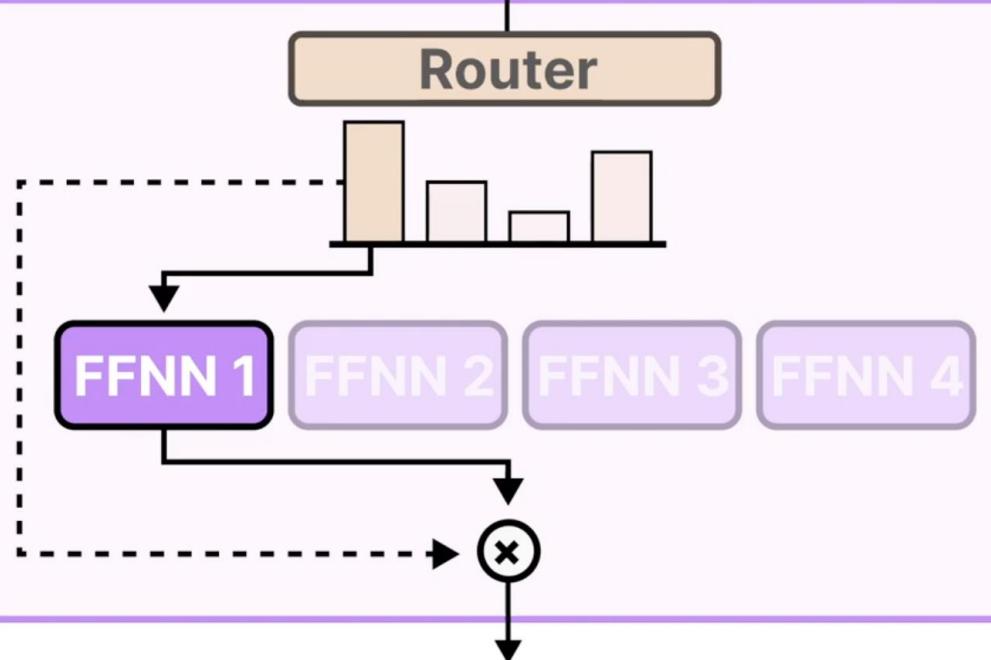
## The Vision Transformer



## The Vision Transformer

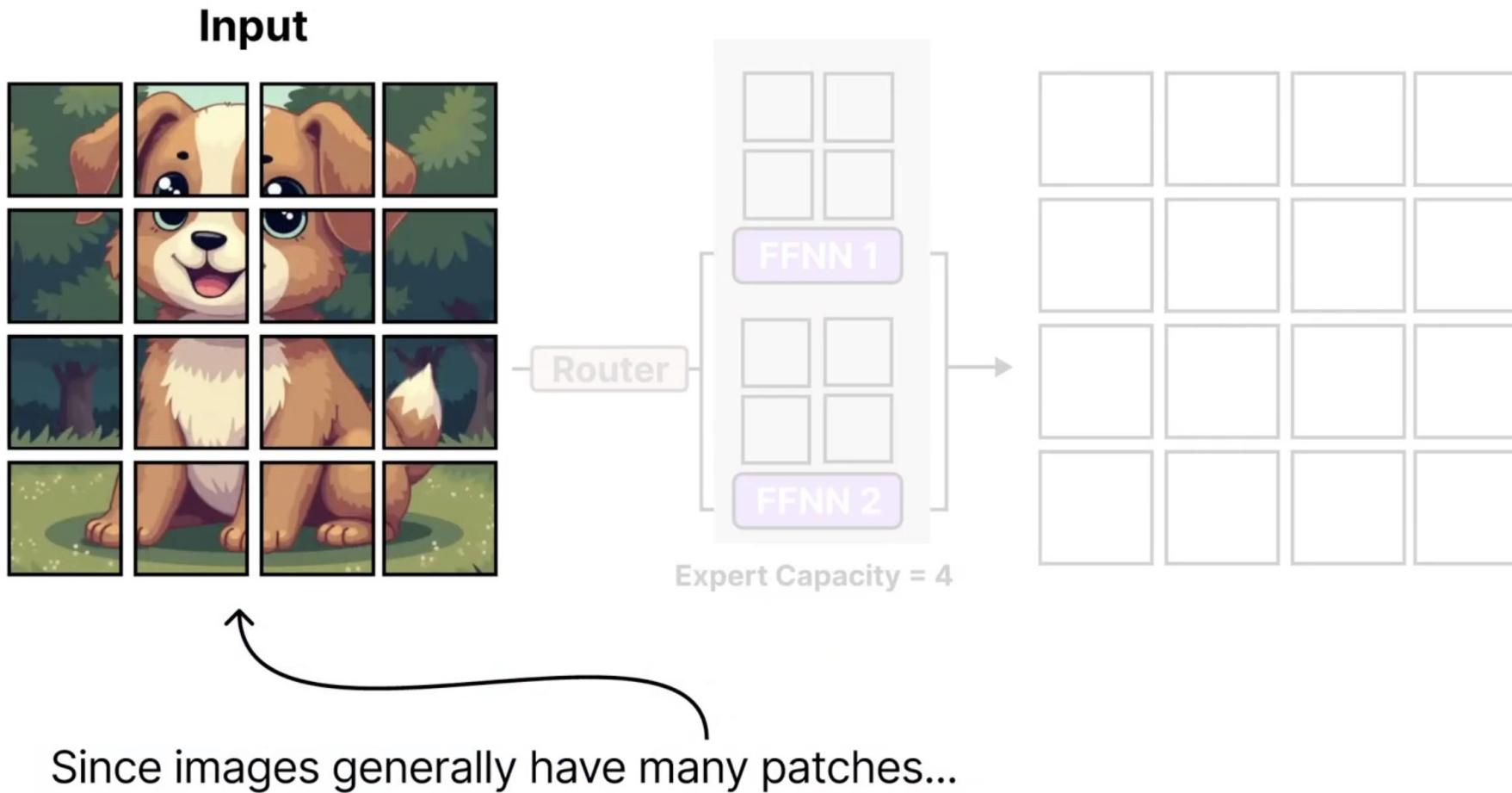


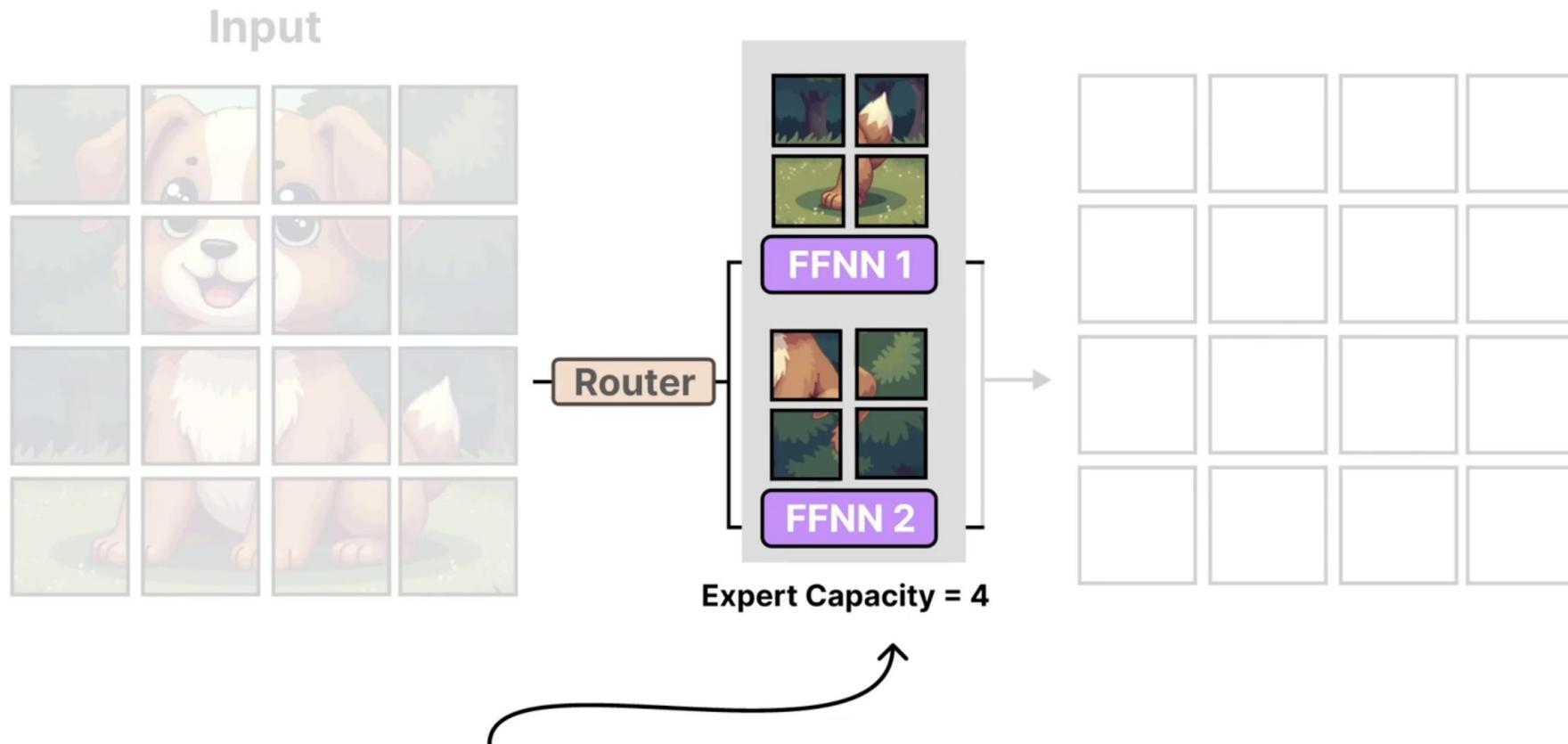
### MoE Layer



Note how we can leverage the existing architecture of Transformers to implement **MoE** in a Vision Model.

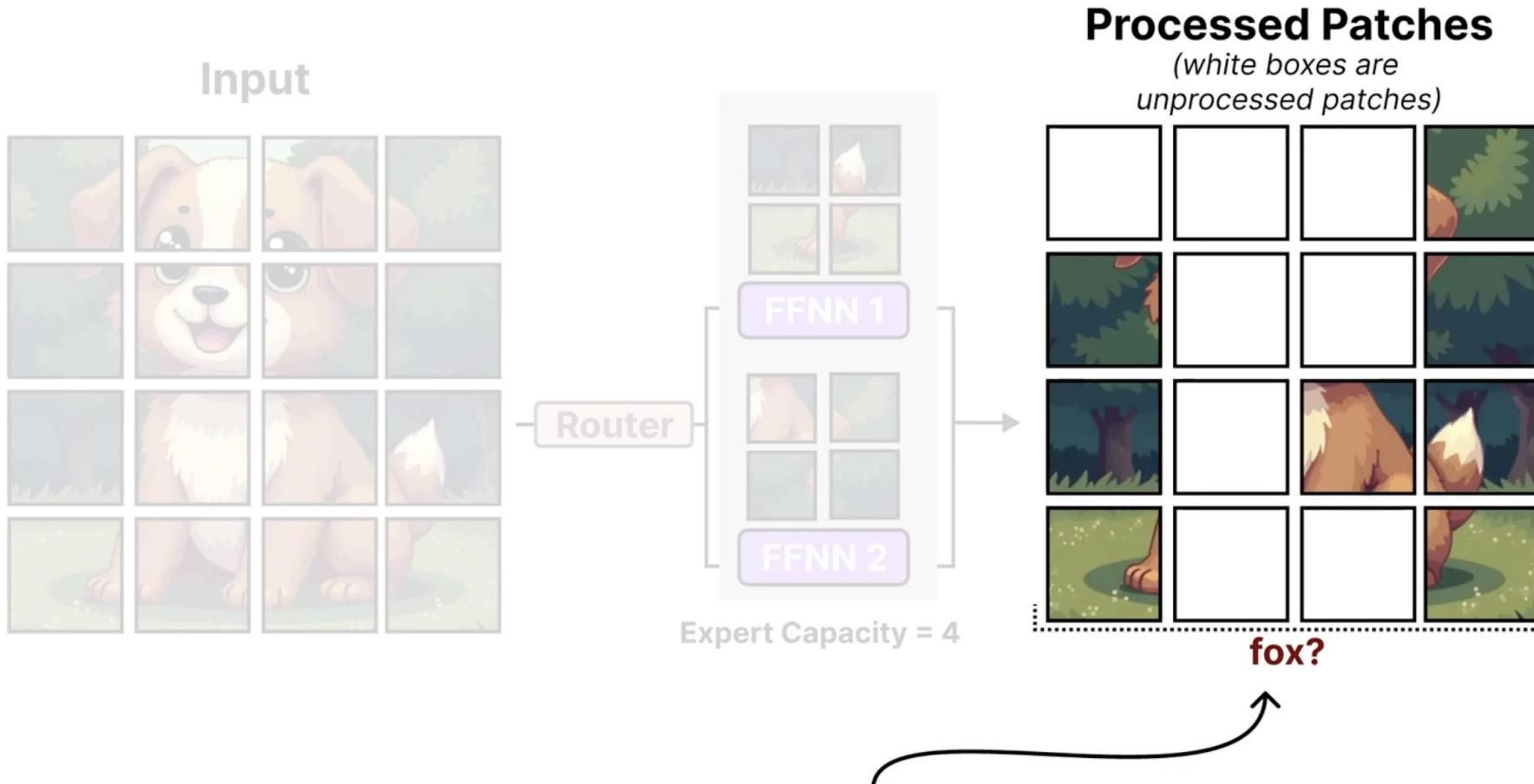






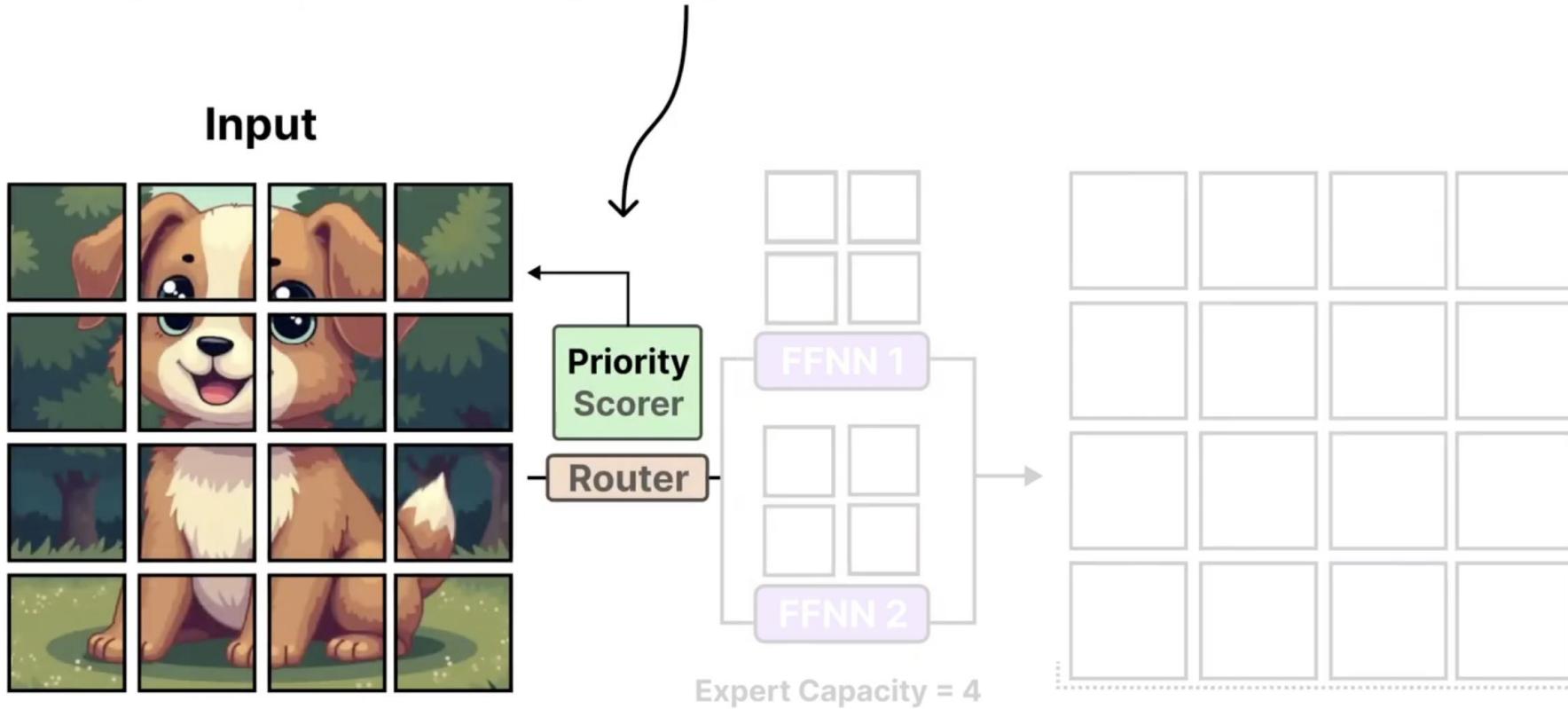
... a low **expert capacity** is used for each expert to reduce hardware constraints.

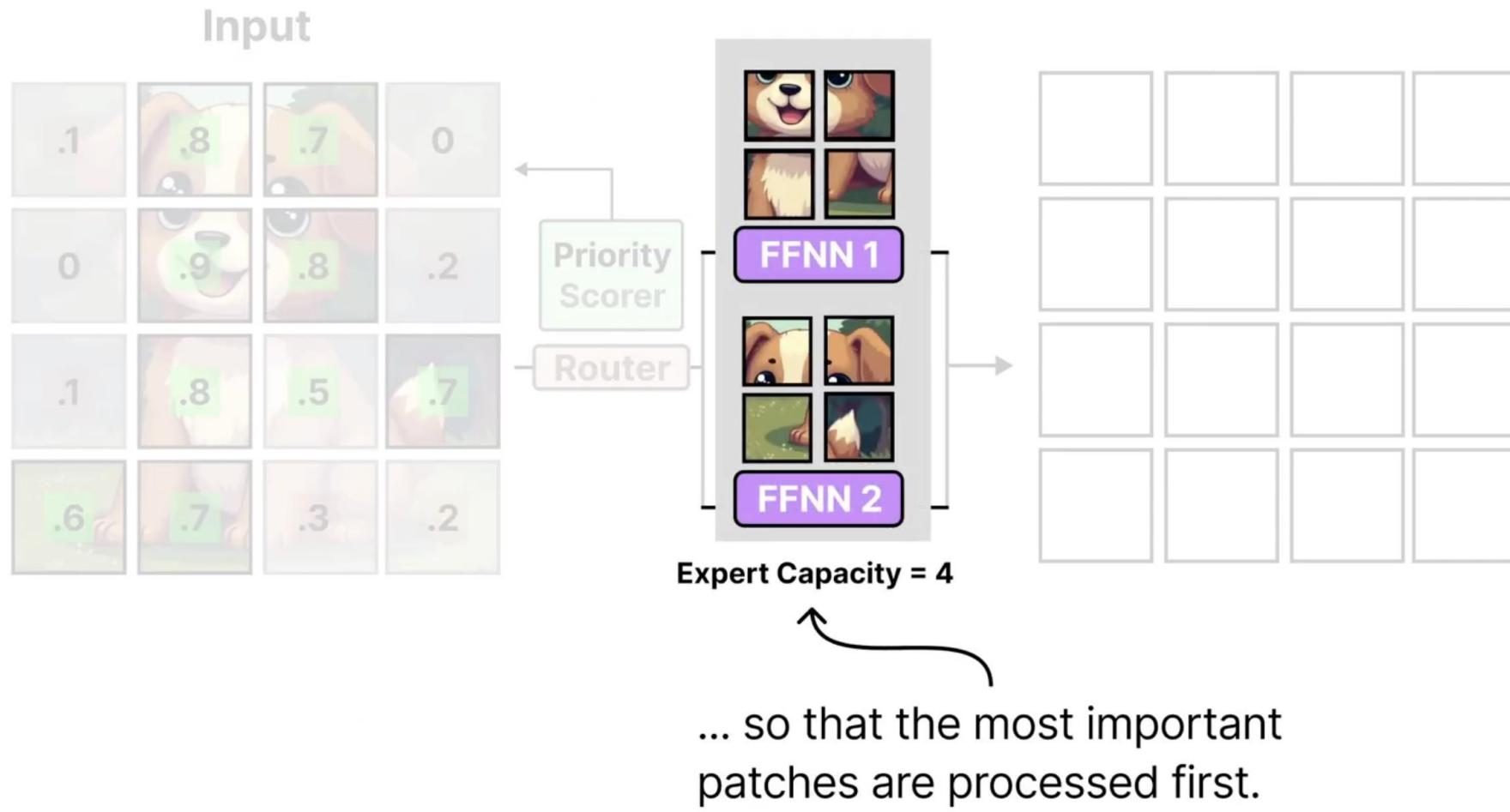


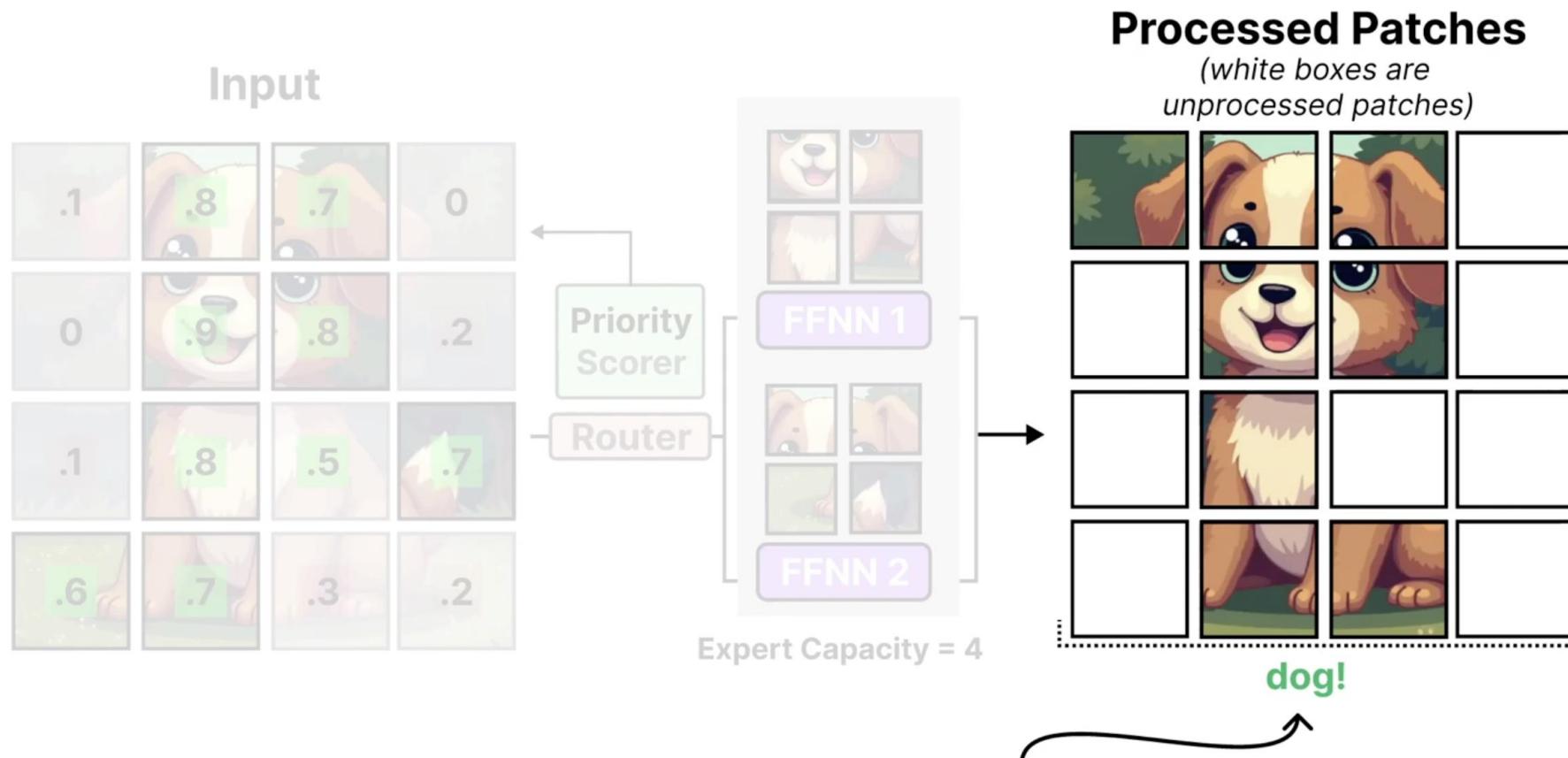


However, a low capacity tends to lead to patches being dropped (akin to **token overflow**).

To keep the capacity low, a **priority scorer**...

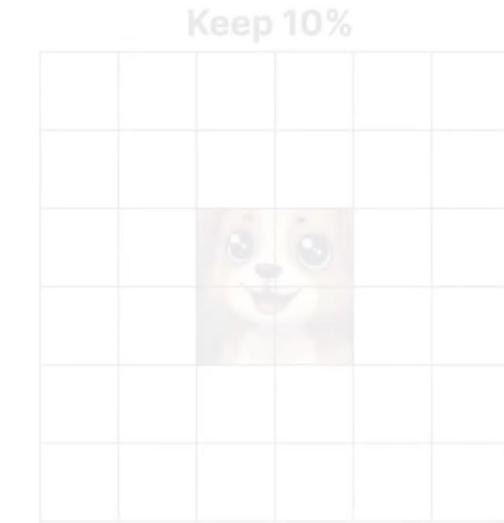
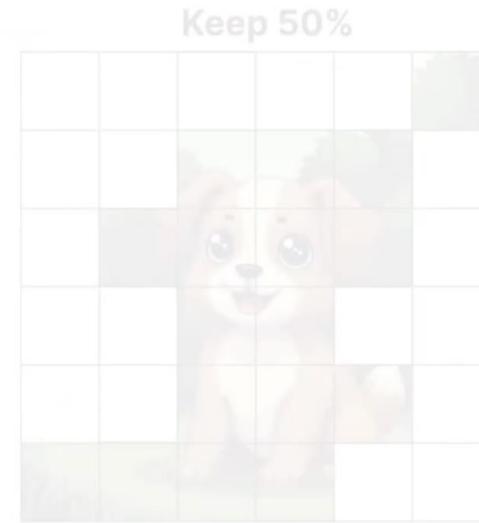






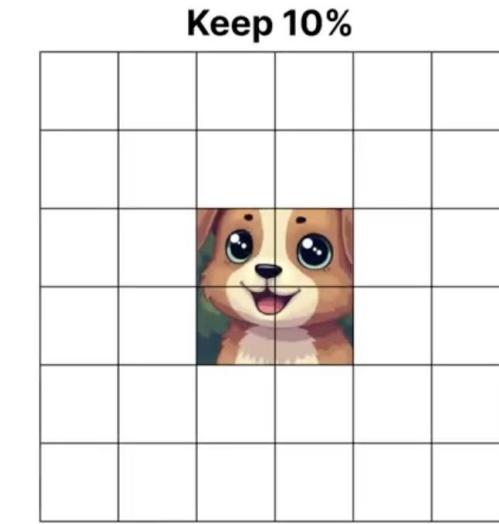
This results in a much more accurate representation of the original image.





As a result, we should still see **important patches** routed if the percentage of tokens decreases.



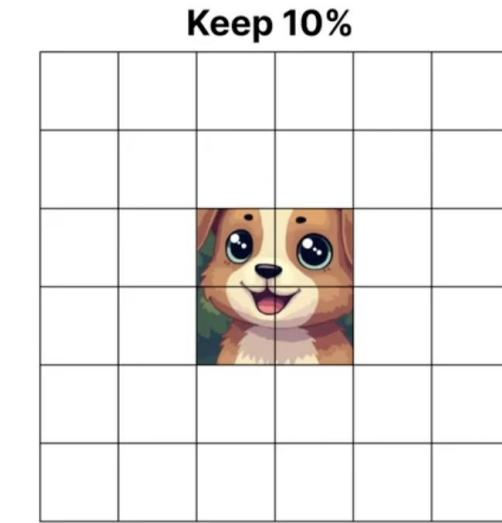
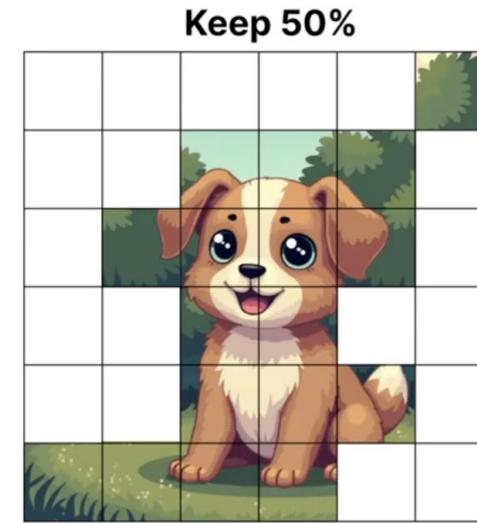


As a result, we should still see **important patches** routed if the percentage of tokens decreases.



02

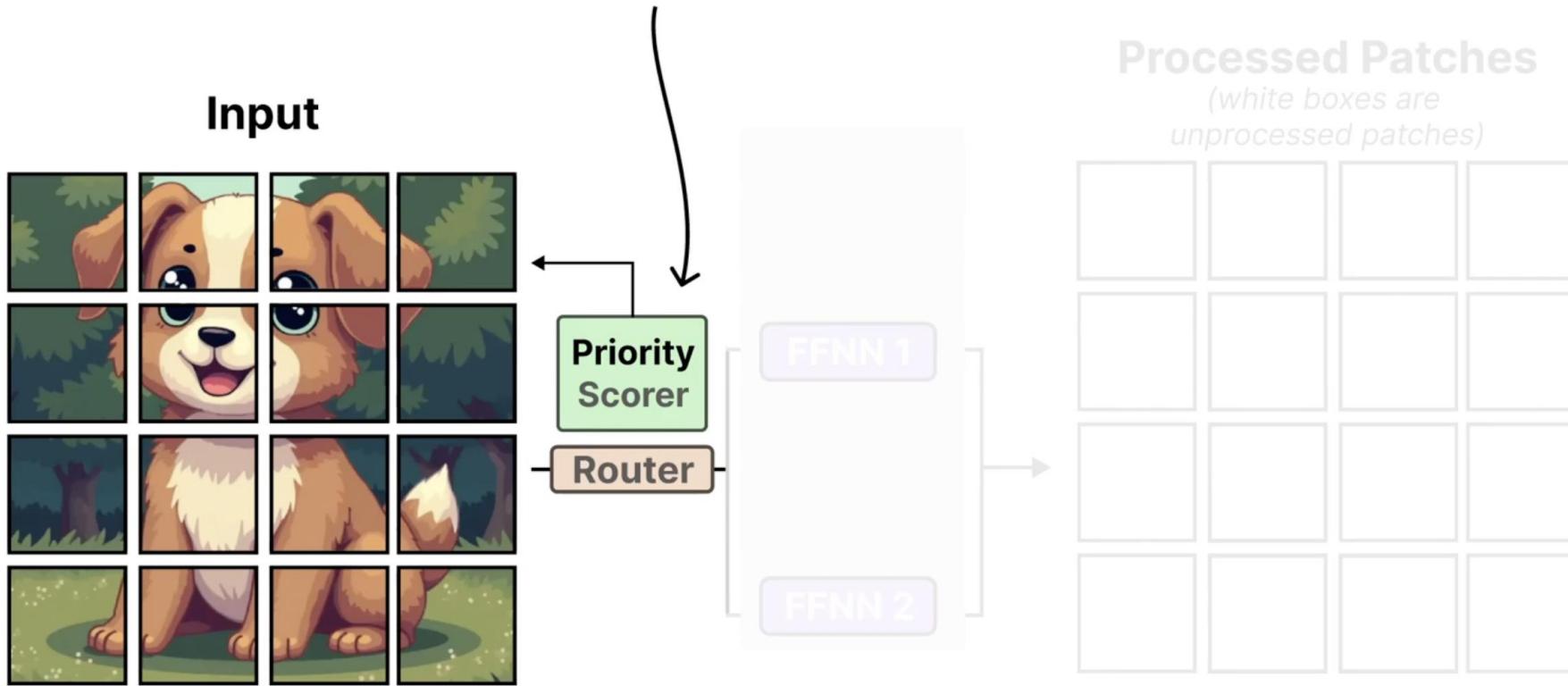
# Soft-MOE 可视化原理

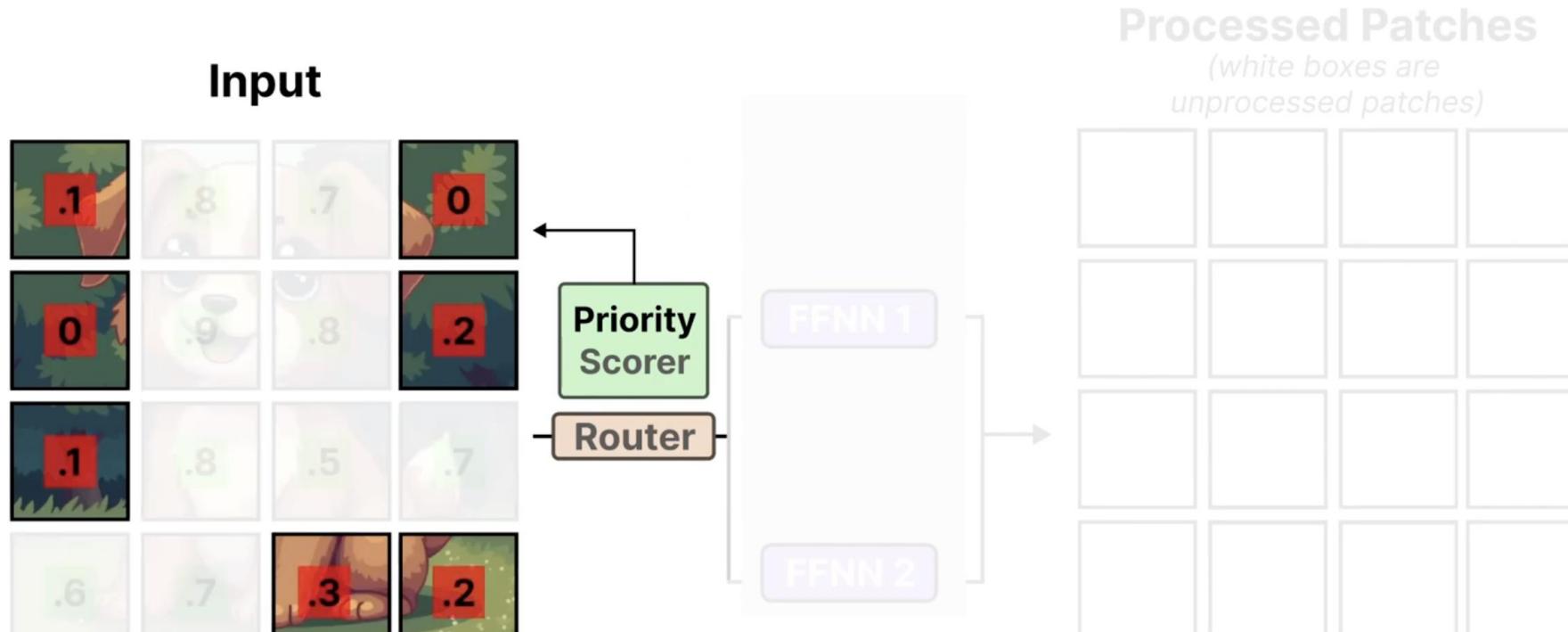


However, **Vision-MoE** still needs to drop tokens, so information is lost. Instead, let's look at an alternative, namely **Soft-MoE**.



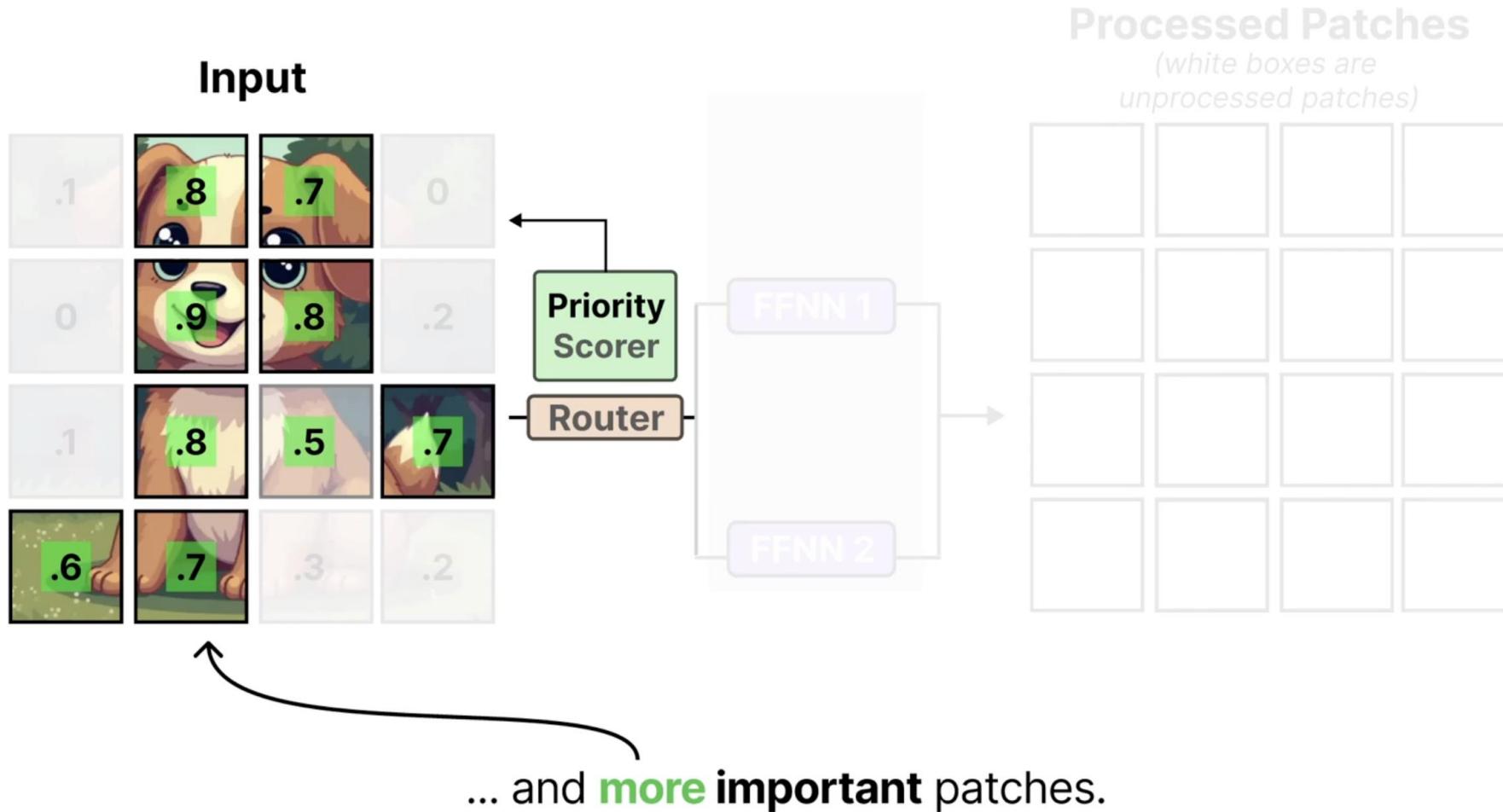
In Vision-MoE, the **priority scorer**...

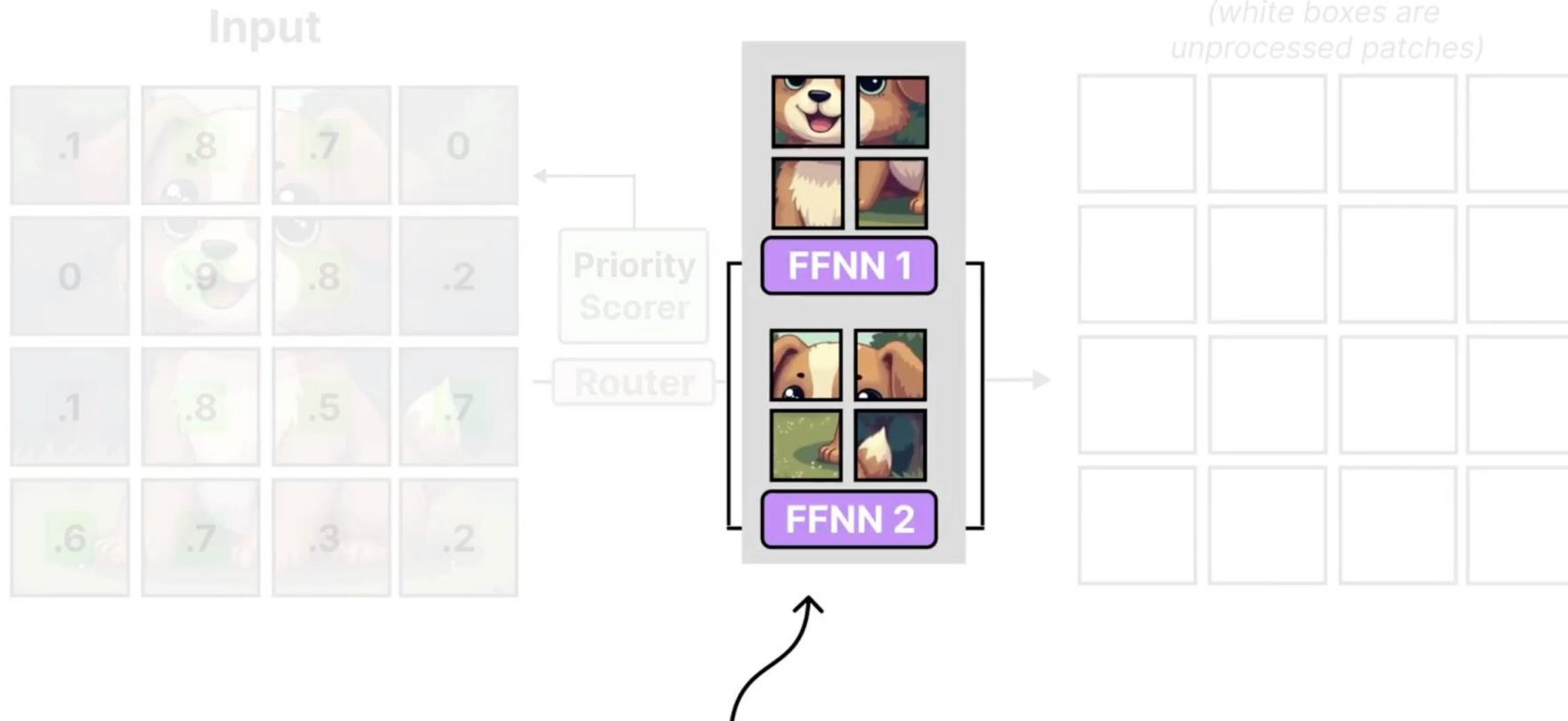




... helps differentiate between **less important** patches...

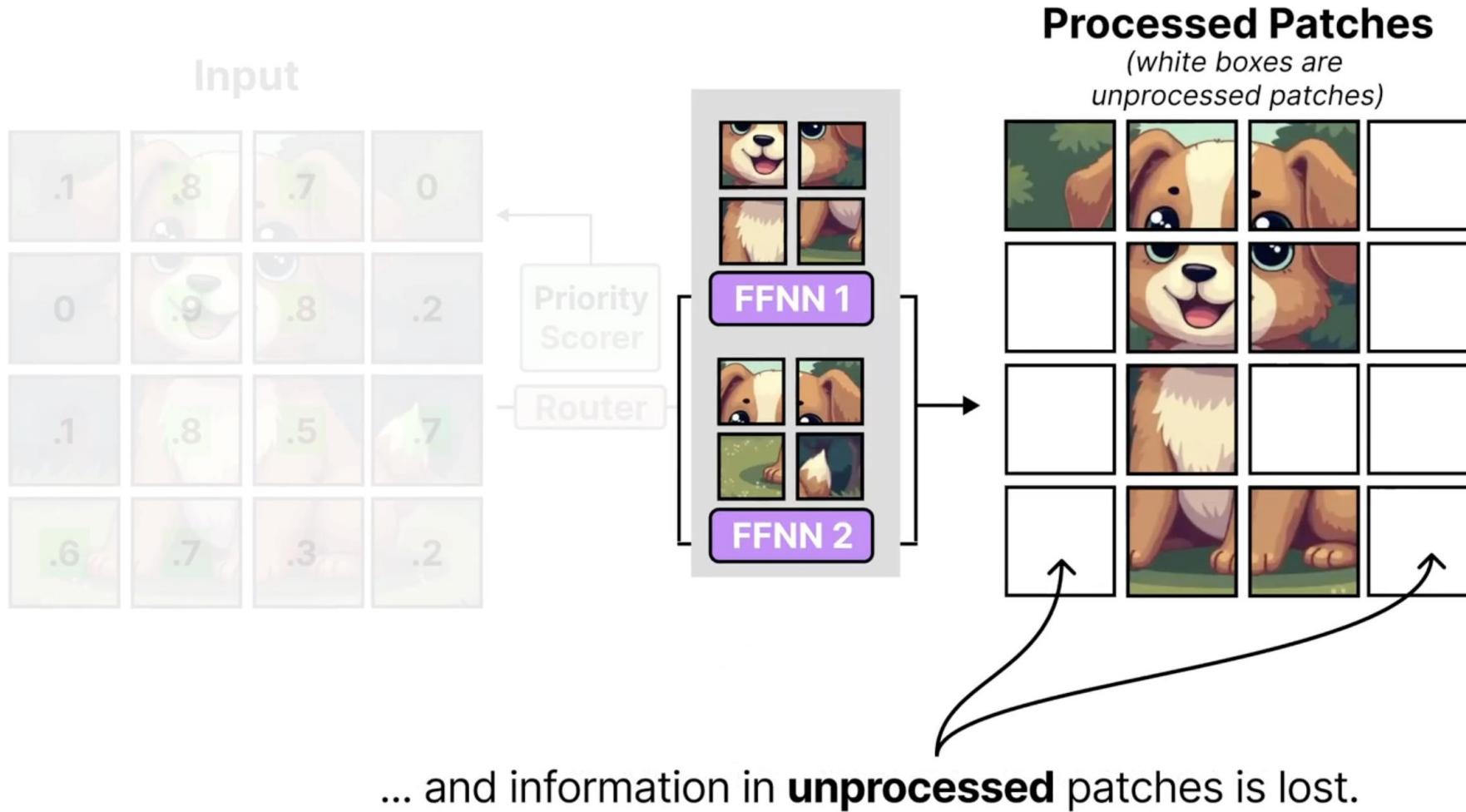


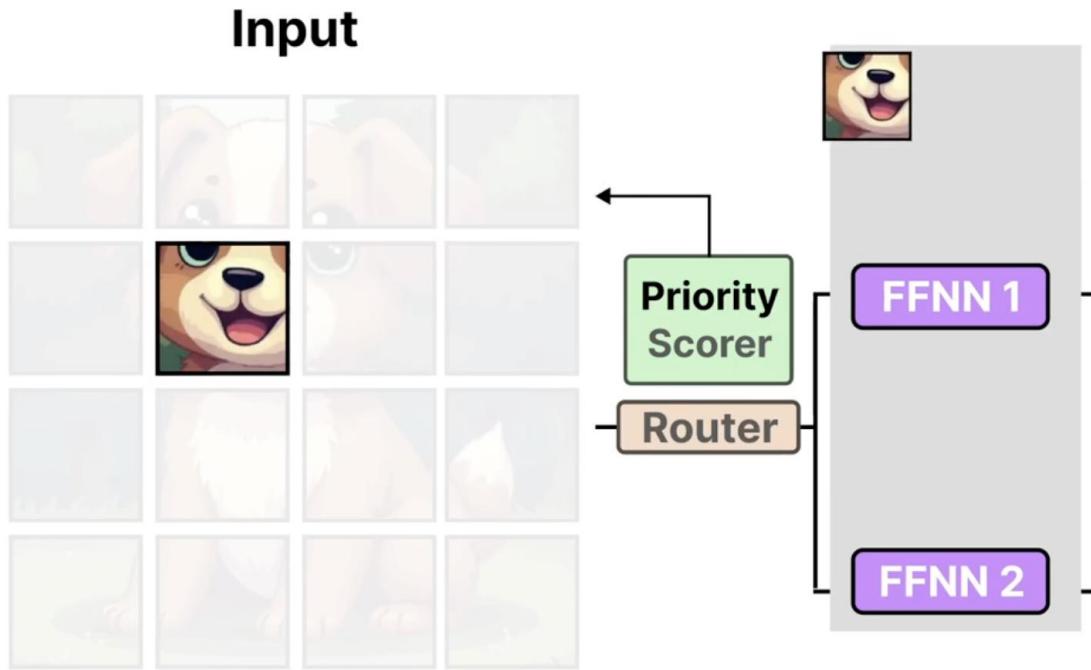




However, a **subset** of patches are assigned to each expert...

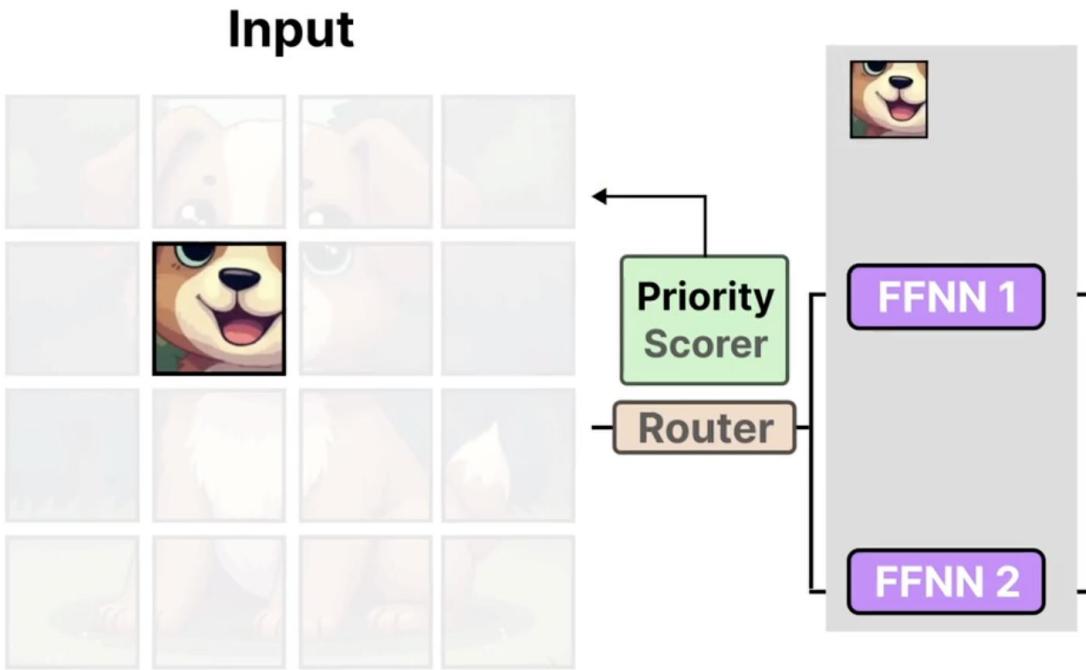






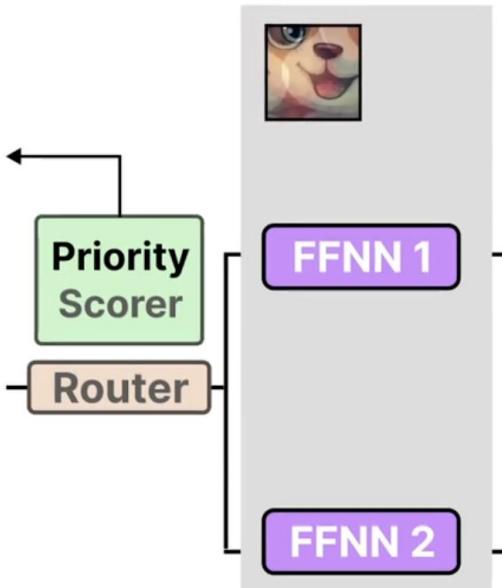
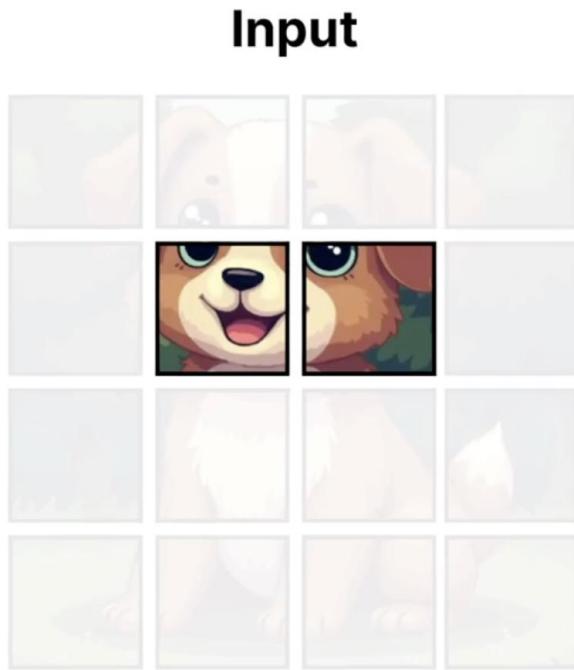
**Soft-Moe** solves this problem, and aims to go from a **discrete** patch assignment...





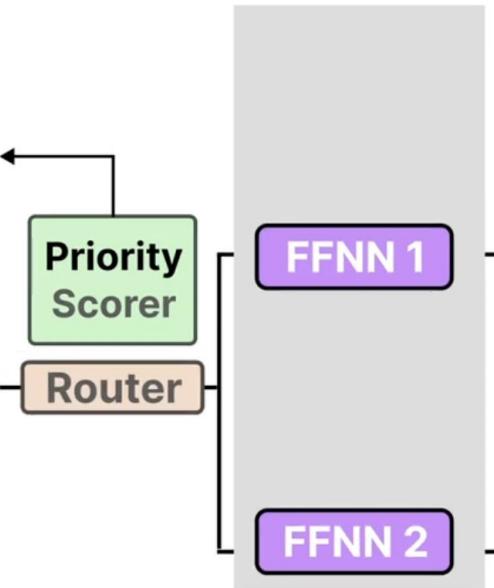
**Soft-Moe** solves this problem, and aims to go from a **discrete** patch assignment...



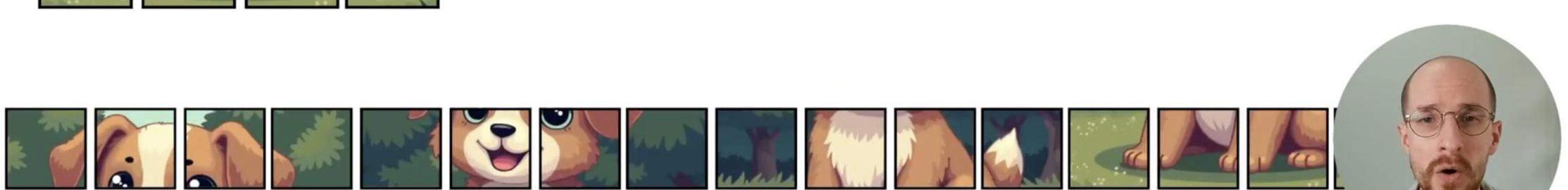


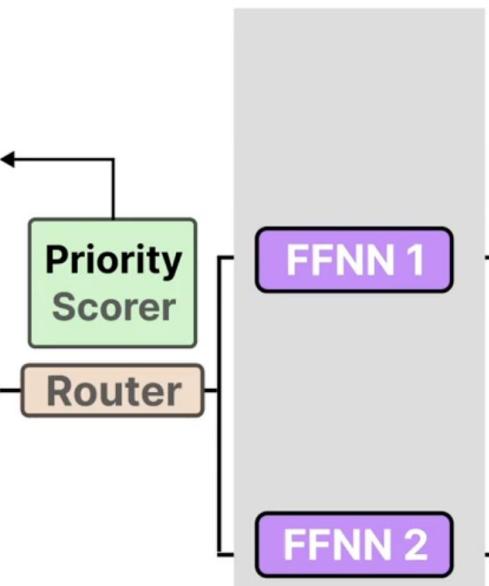
... to a **soft** patch assignment by mixing patches.



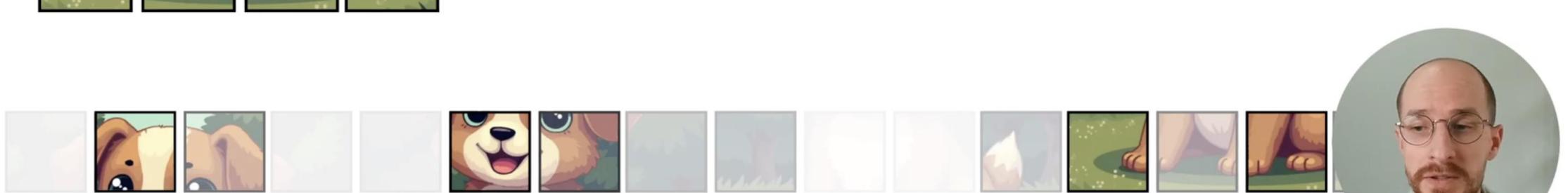


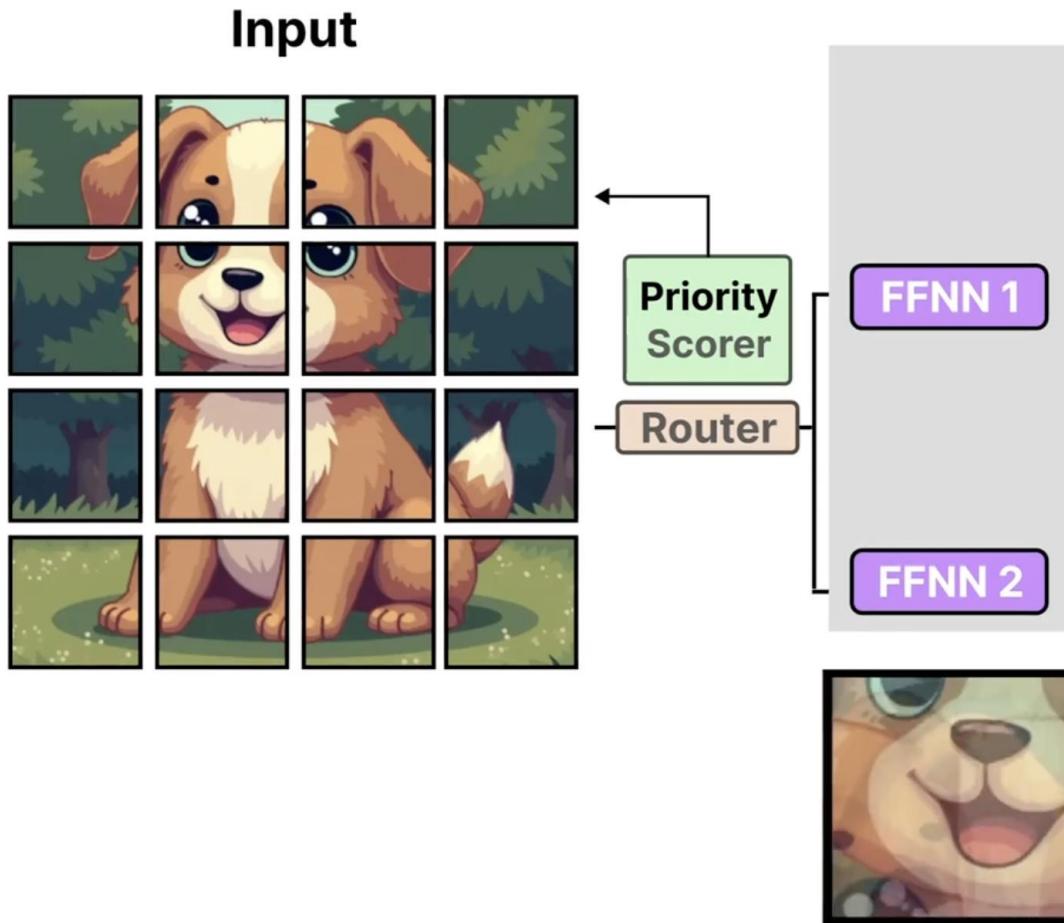
In practice, **all patches** are used instead of only a subset to create this soft patch.





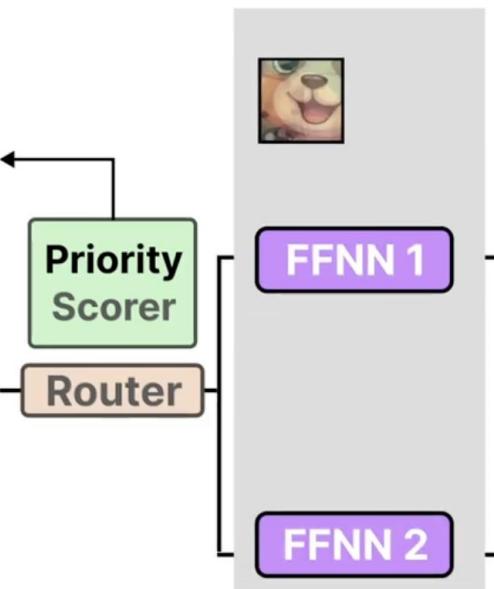
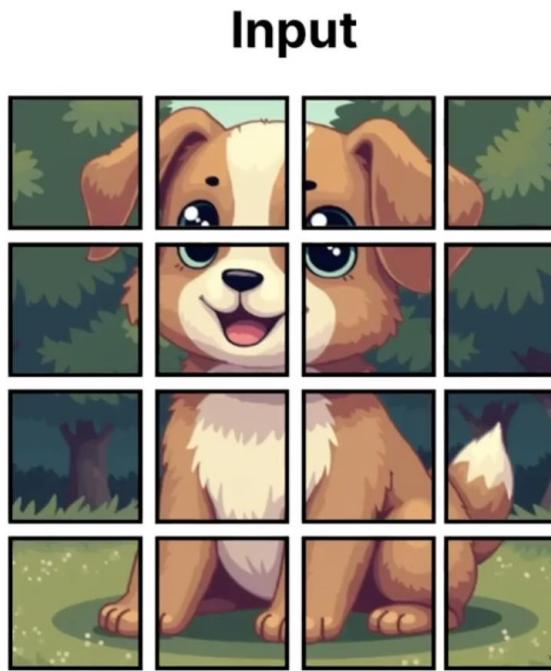
We **weigh** each patch using the priority scorer...





... and then take a **linear combination** of these weighted patches...

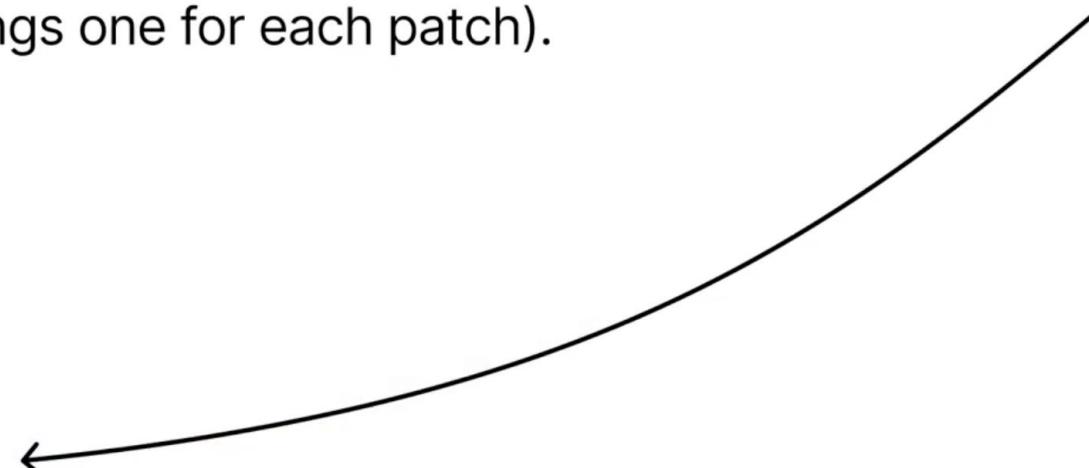
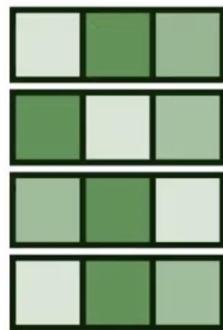




... before we send the **soft patch** to the chosen expert.



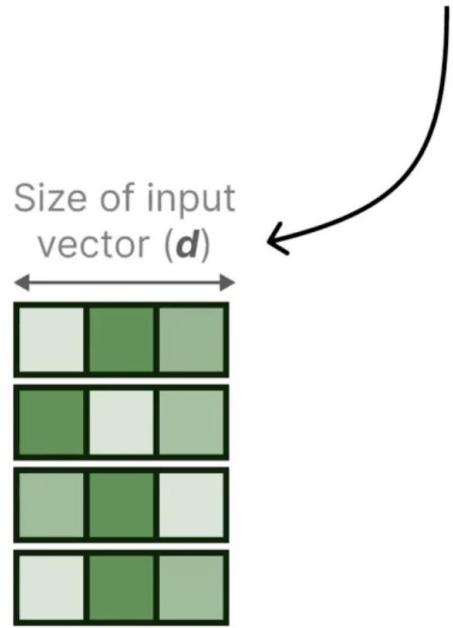
To create these soft patches, **Soft-MoE** takes the input **X** (patch embeddings one for each patch).



**X**  
(patch **embeddings**)



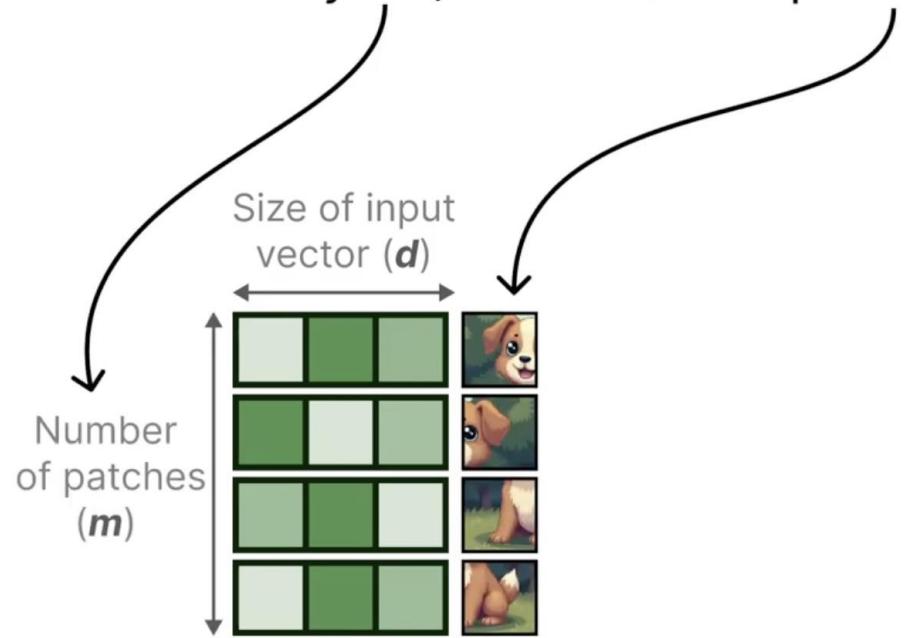
With dimensions  $d$  (size of input vector)...



X  
(patch **embeddings**)



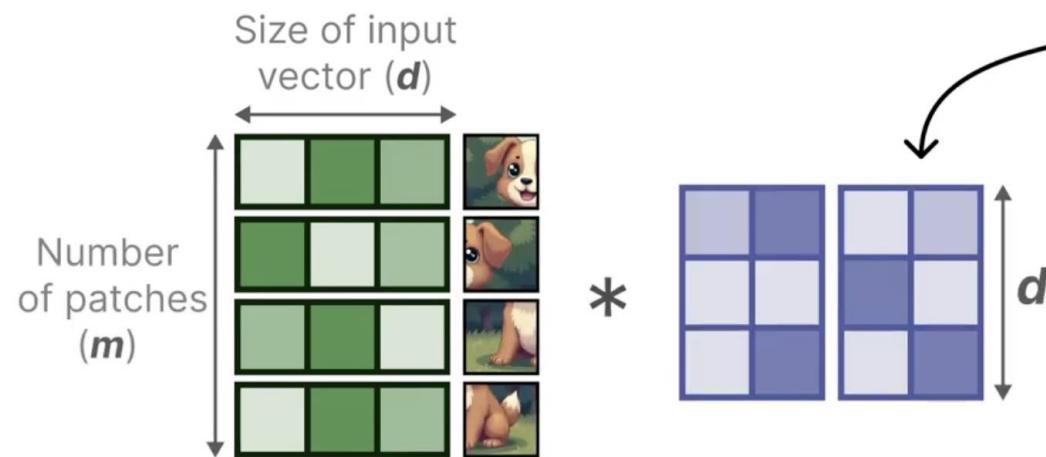
.. by  $m$  (the number of patches).



X  
(patch **embeddings**)



**X** is then multiplied by a learnable matrix  $\Phi$  with dimensions  $d$  by...

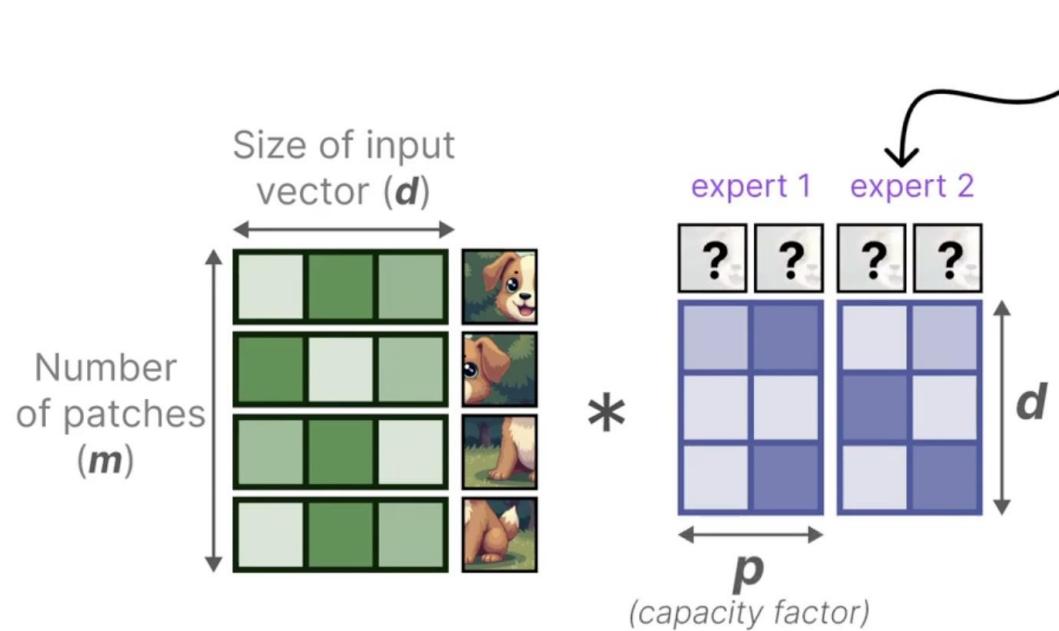


**X**  
(patch **embeddings**)

**$\Phi$**   
(learnable  
matrix)



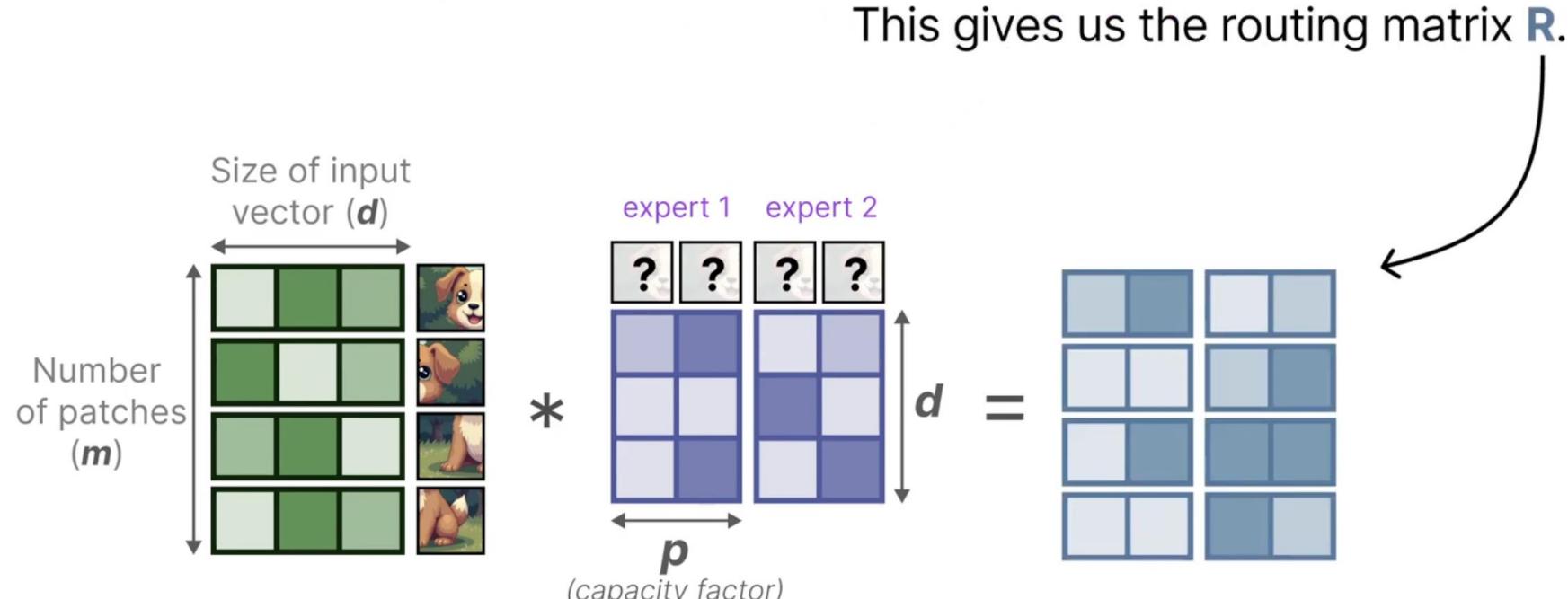
... the capacity factor for each **expert**.



**X**  
(patch **embeddings**)

**$\Phi$**   
(learnable  
matrix)





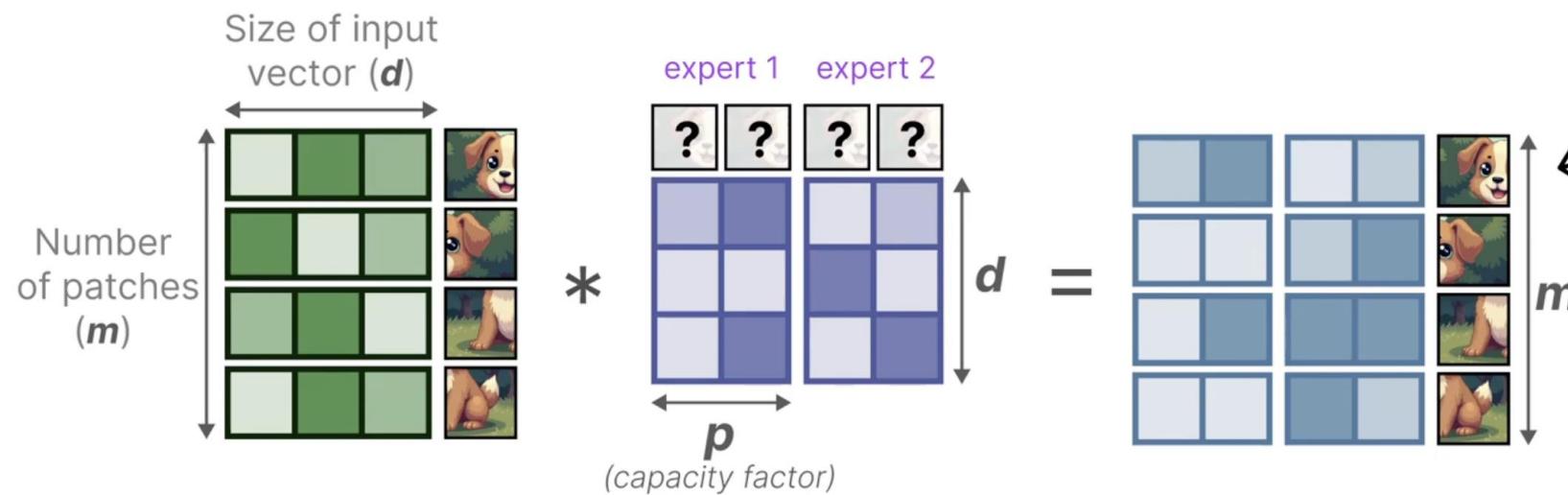
**X**  
(patch **embeddings**)

**$\Phi$**   
(learnable  
matrix)

**R**  
(routing matrix)



... which tells us how related a certain **patch** is...

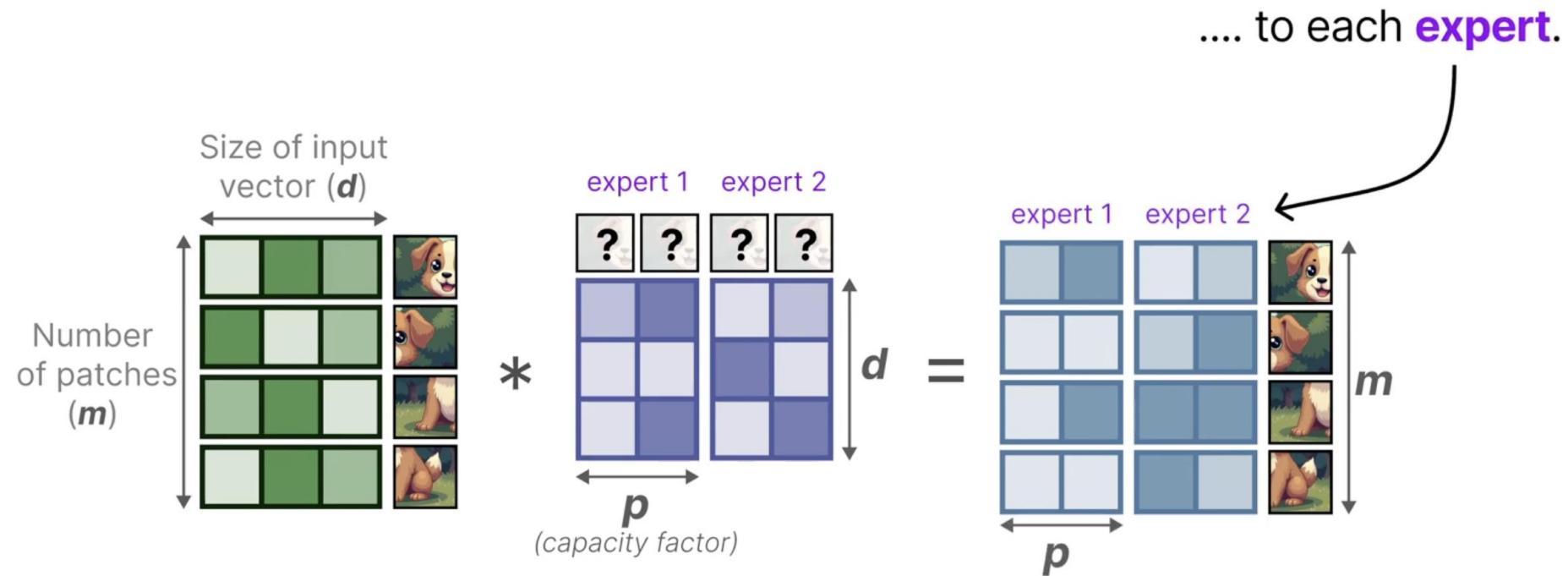


**X**  
(patch **embeddings**)

**$\Phi$**   
(learnable  
matrix)

**R**  
(routing matrix)





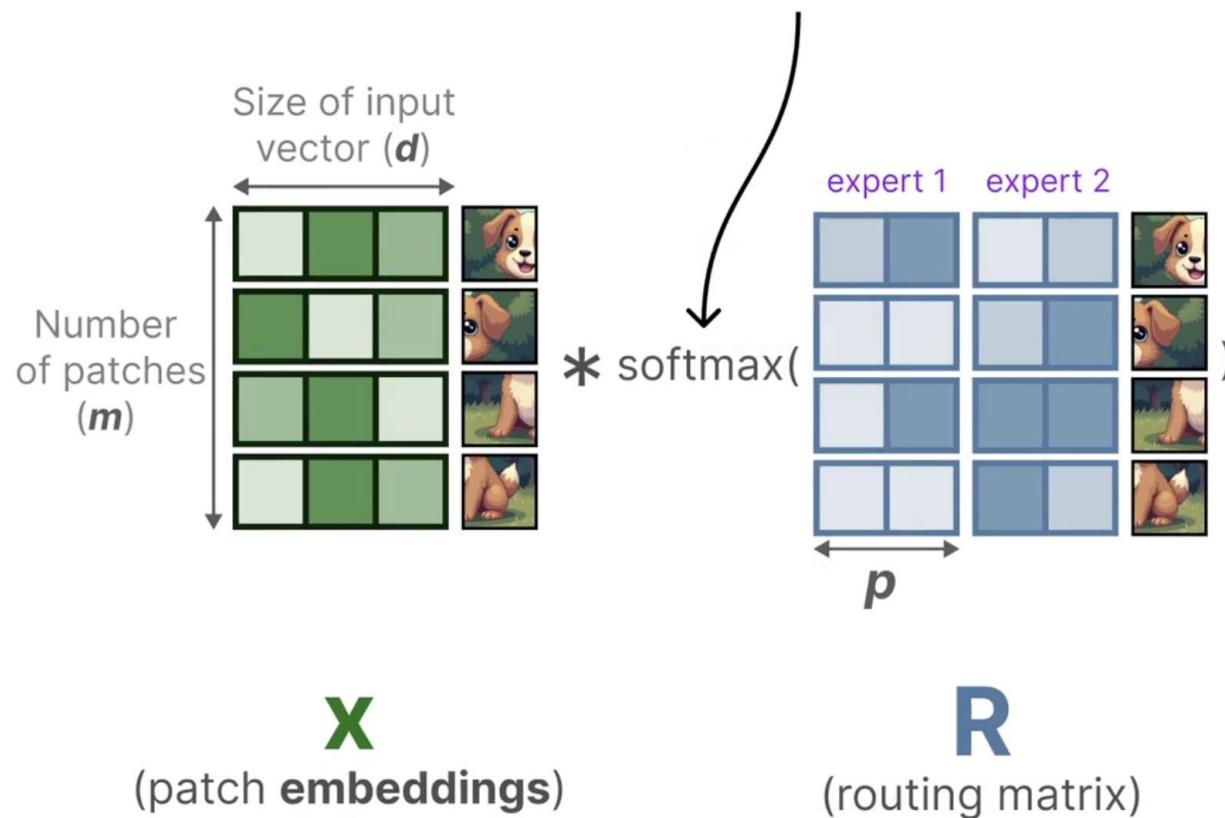
**X**  
(patch **embeddings**)

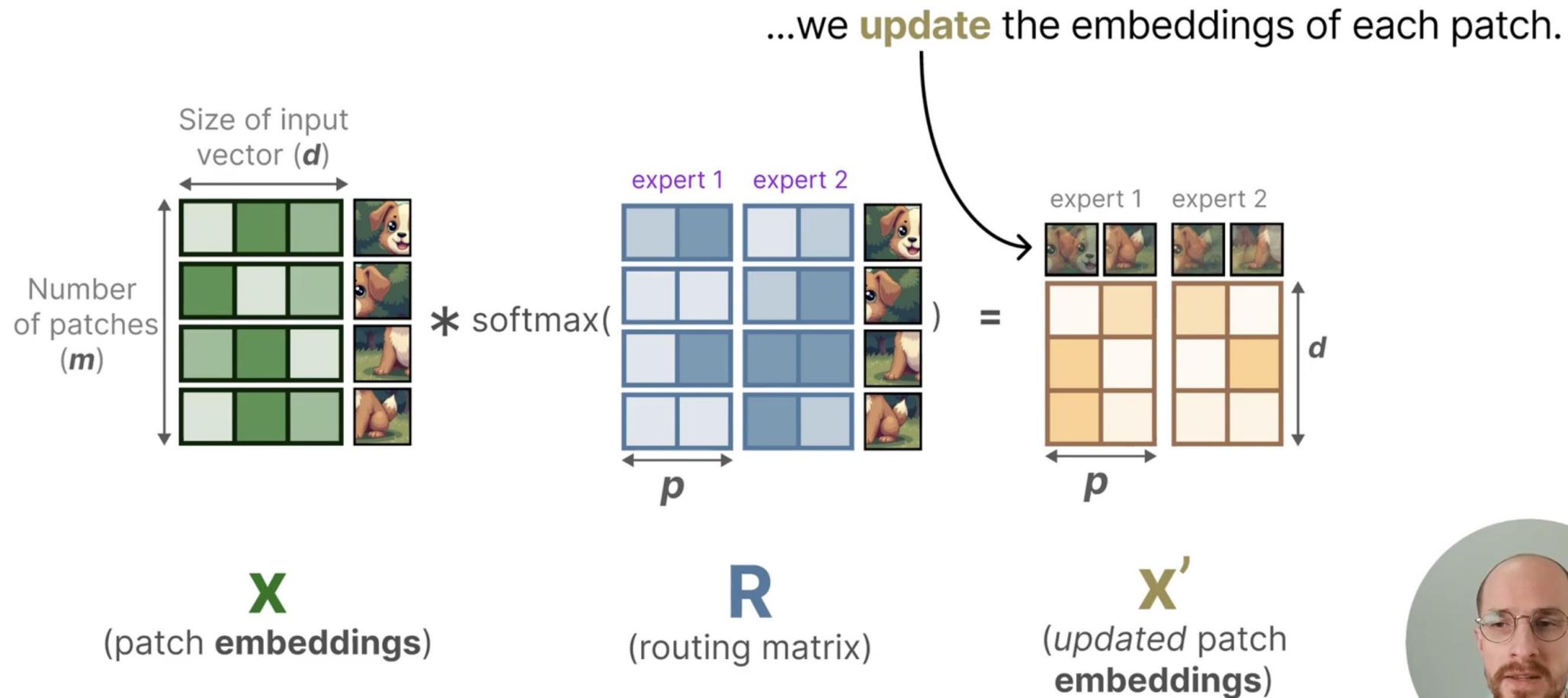
**$\Phi$**   
(learnable  
matrix)

**R**  
(routing matrix)

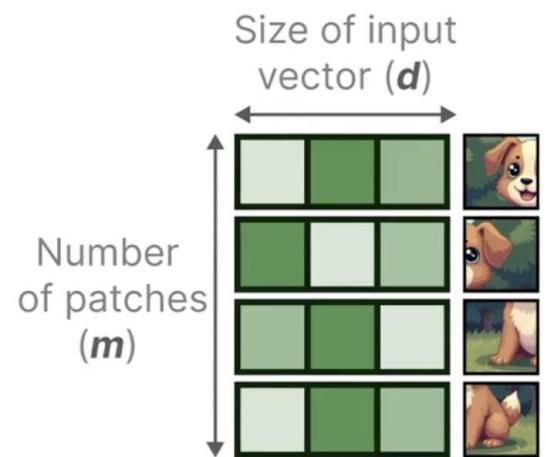


By then taking the **softmax** of the **router matrix** (on the columns) ....

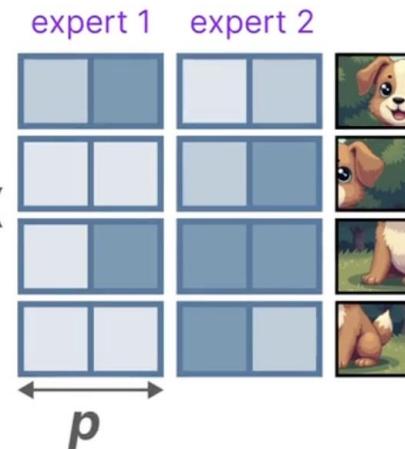




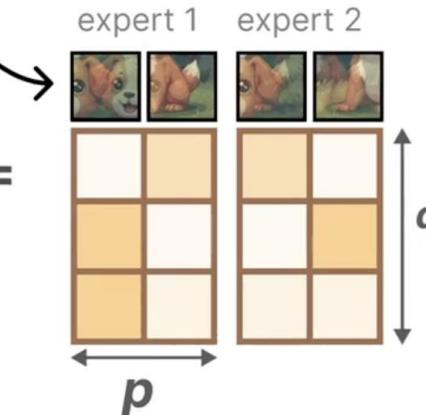
As a result, they are a **linear combination** of the input patches as we saw before.



\* softmax(



)



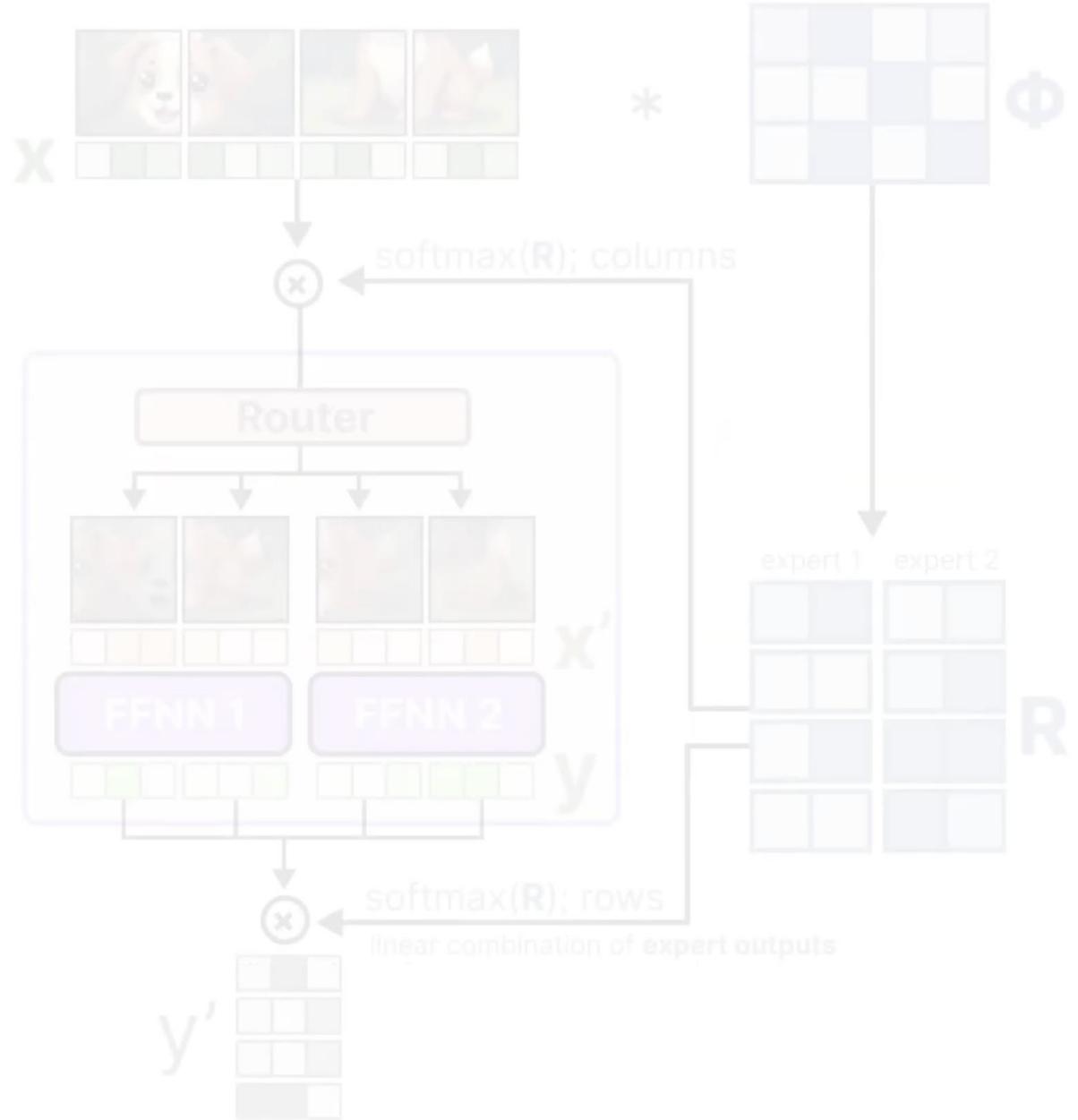
**X**  
(patch **embeddings**)

**R**  
(routing matrix)

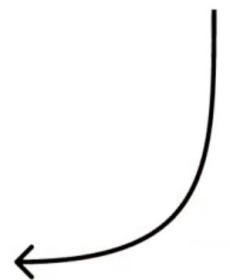
**X'**  
(updated patch  
embeddings)



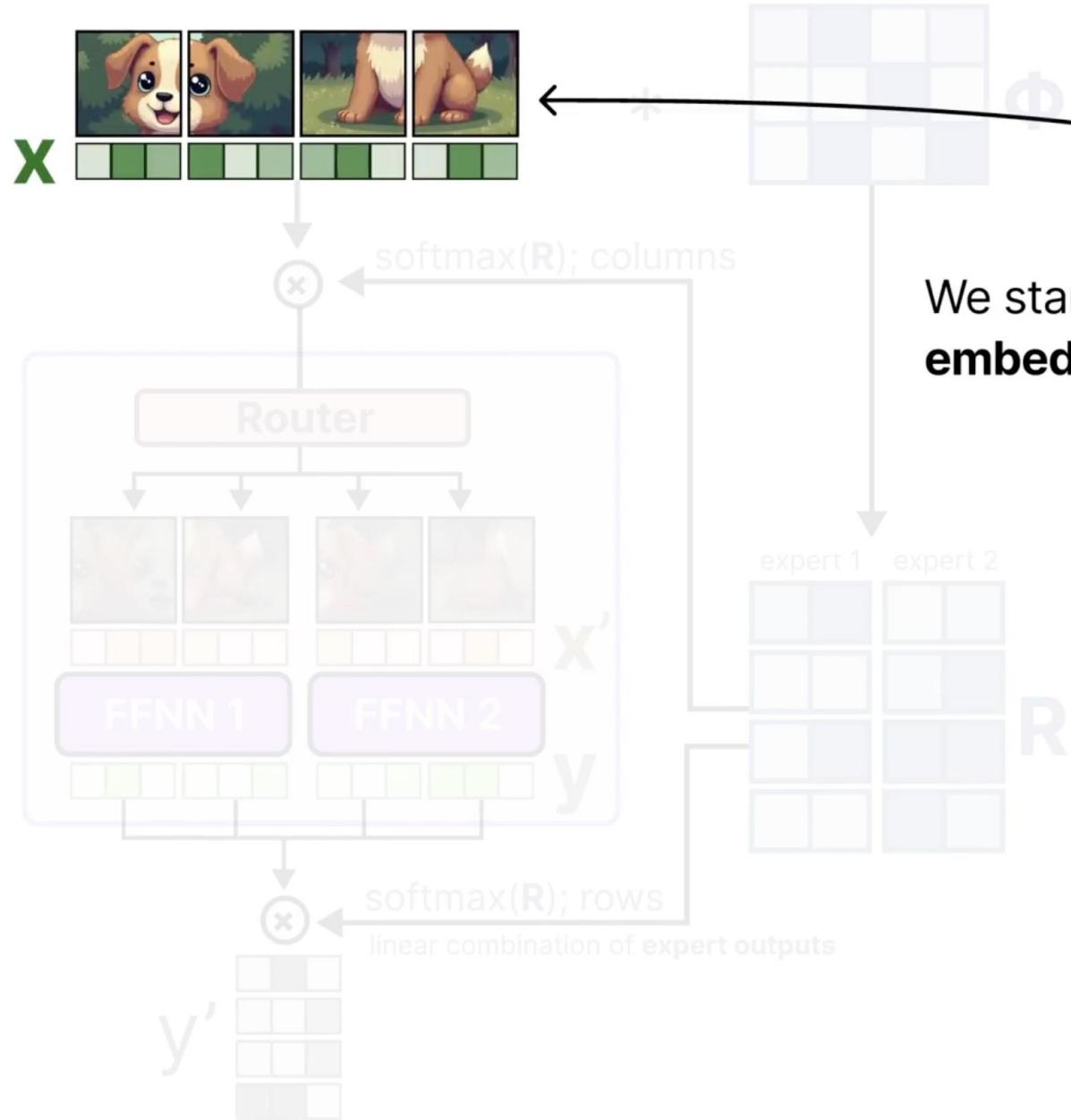
## From Sparse to Soft MoE

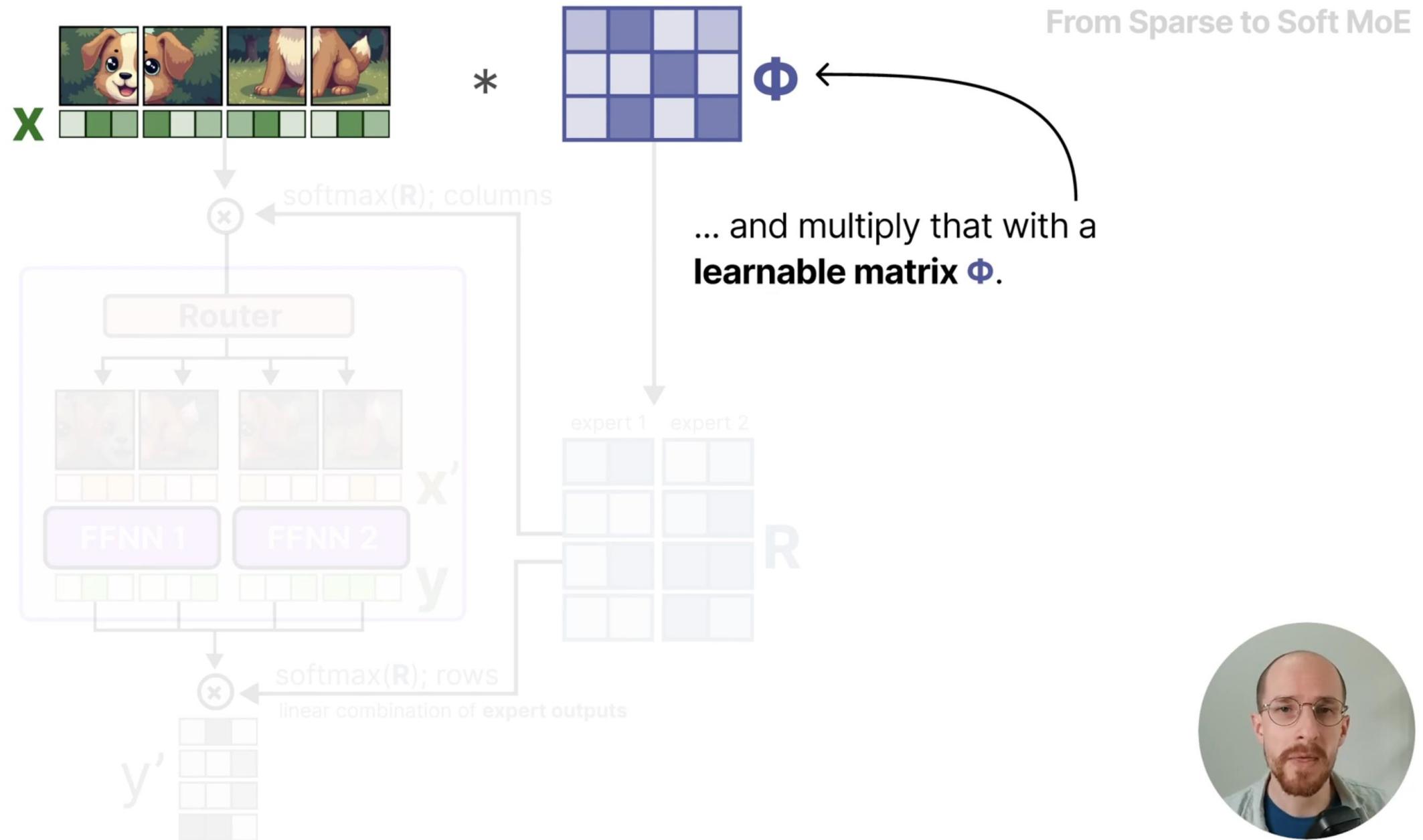


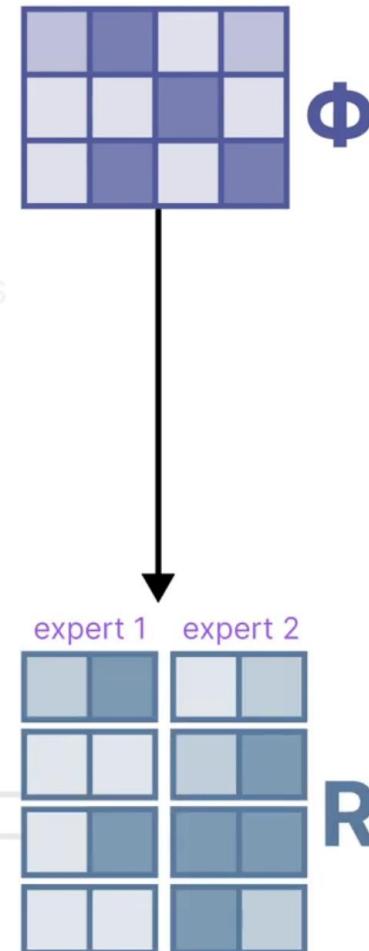
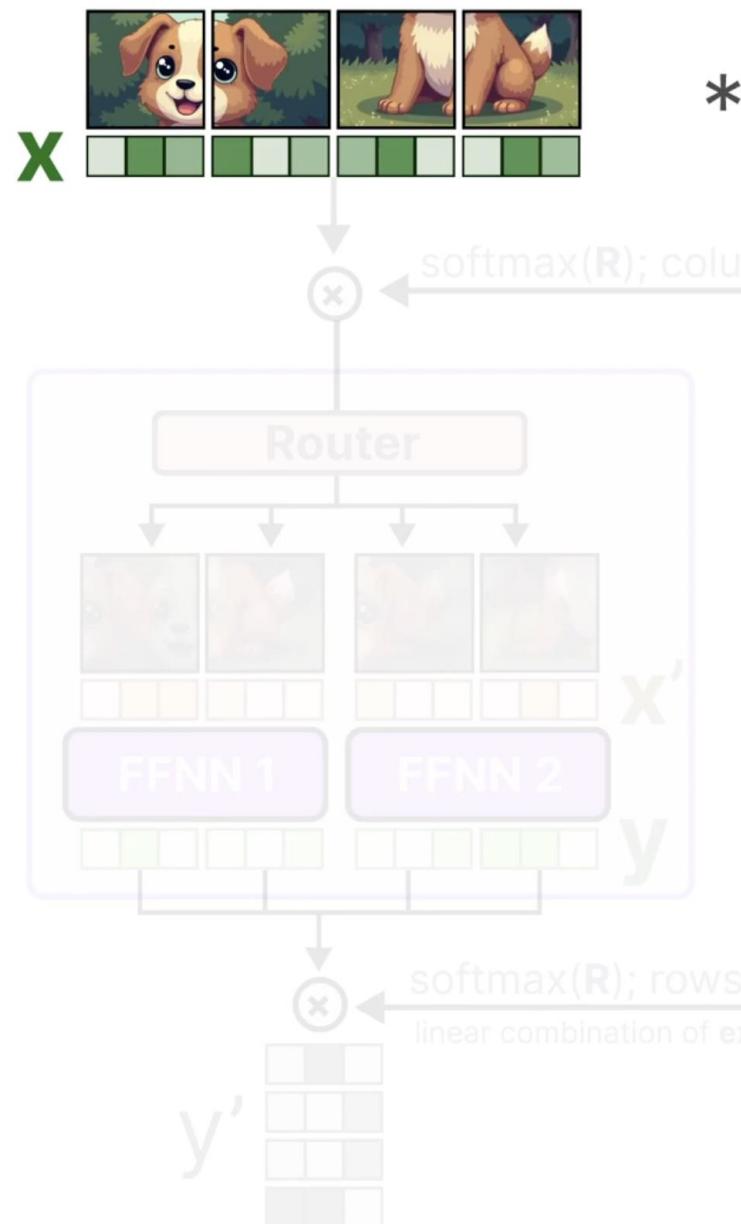
Let's go through this  
**step-by-step.**



## From Sparse to Soft MoE





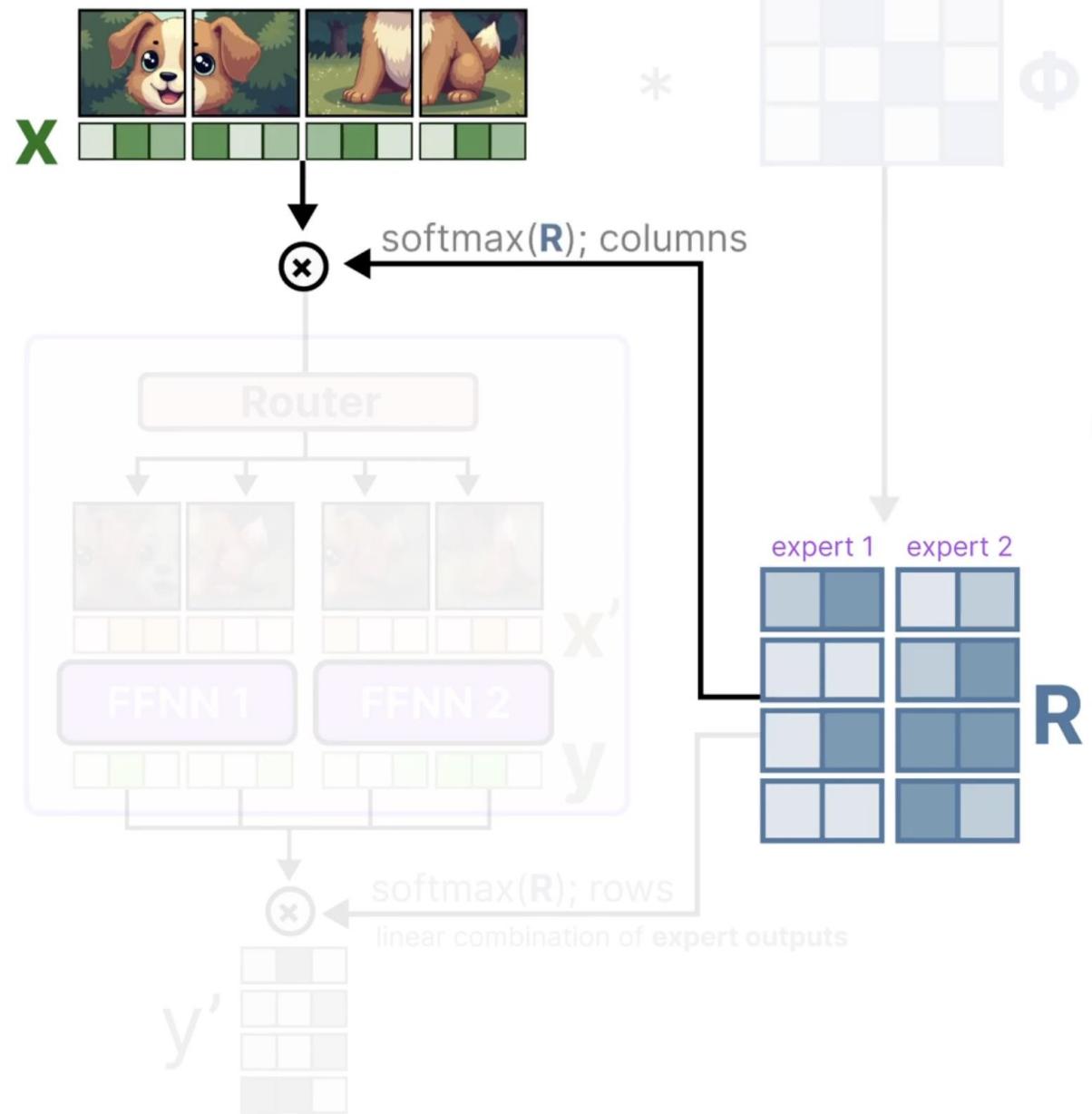


From Sparse to Soft MoE

This gives us the **router matrix  $R$** , which tells us how related a certain token is to a given expert.



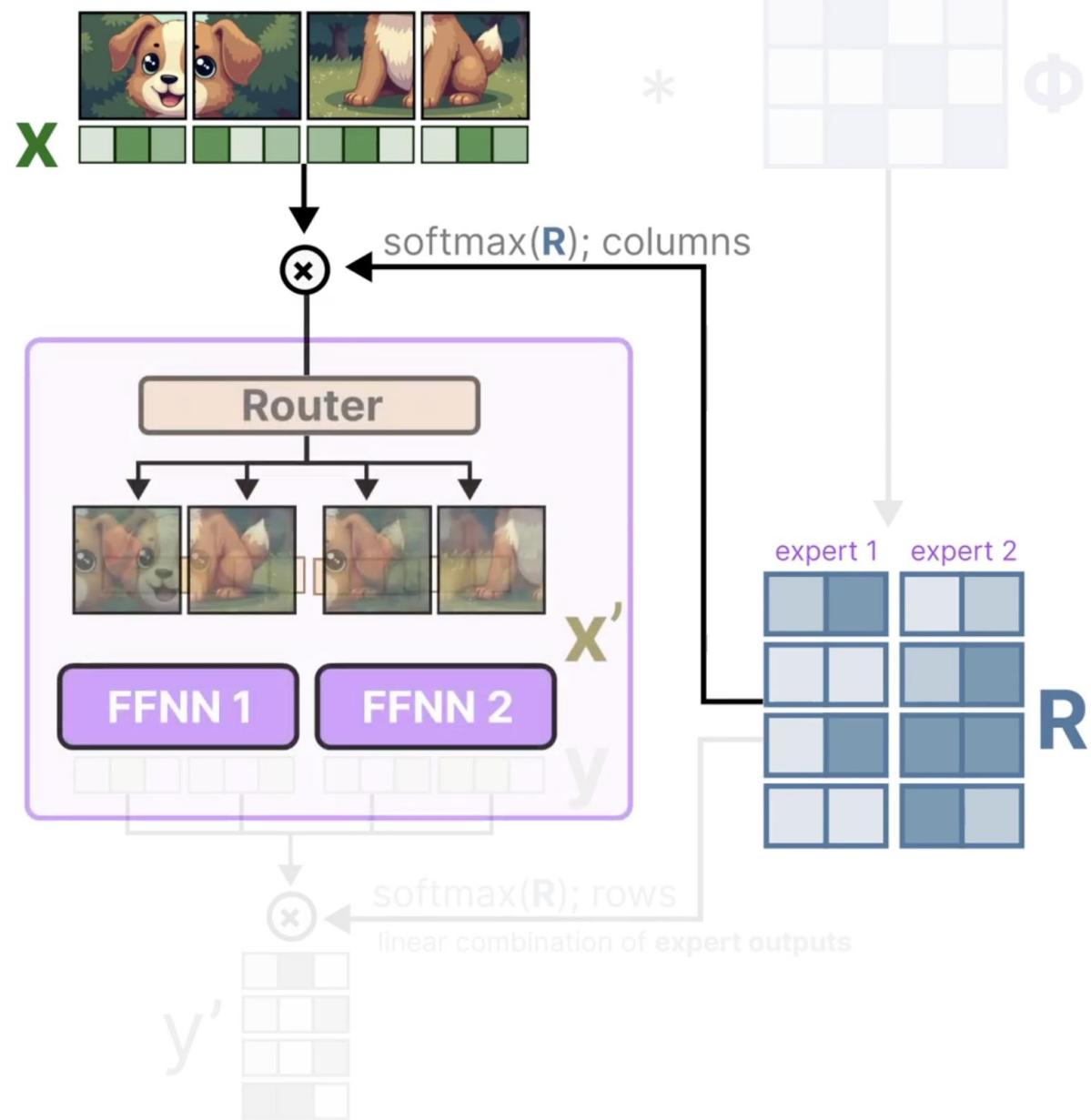
## From Sparse to Soft MoE



Taking the softmax of  $\mathbf{R}$  and multiplying that with the input gives us a **linear combination of patch inputs**.



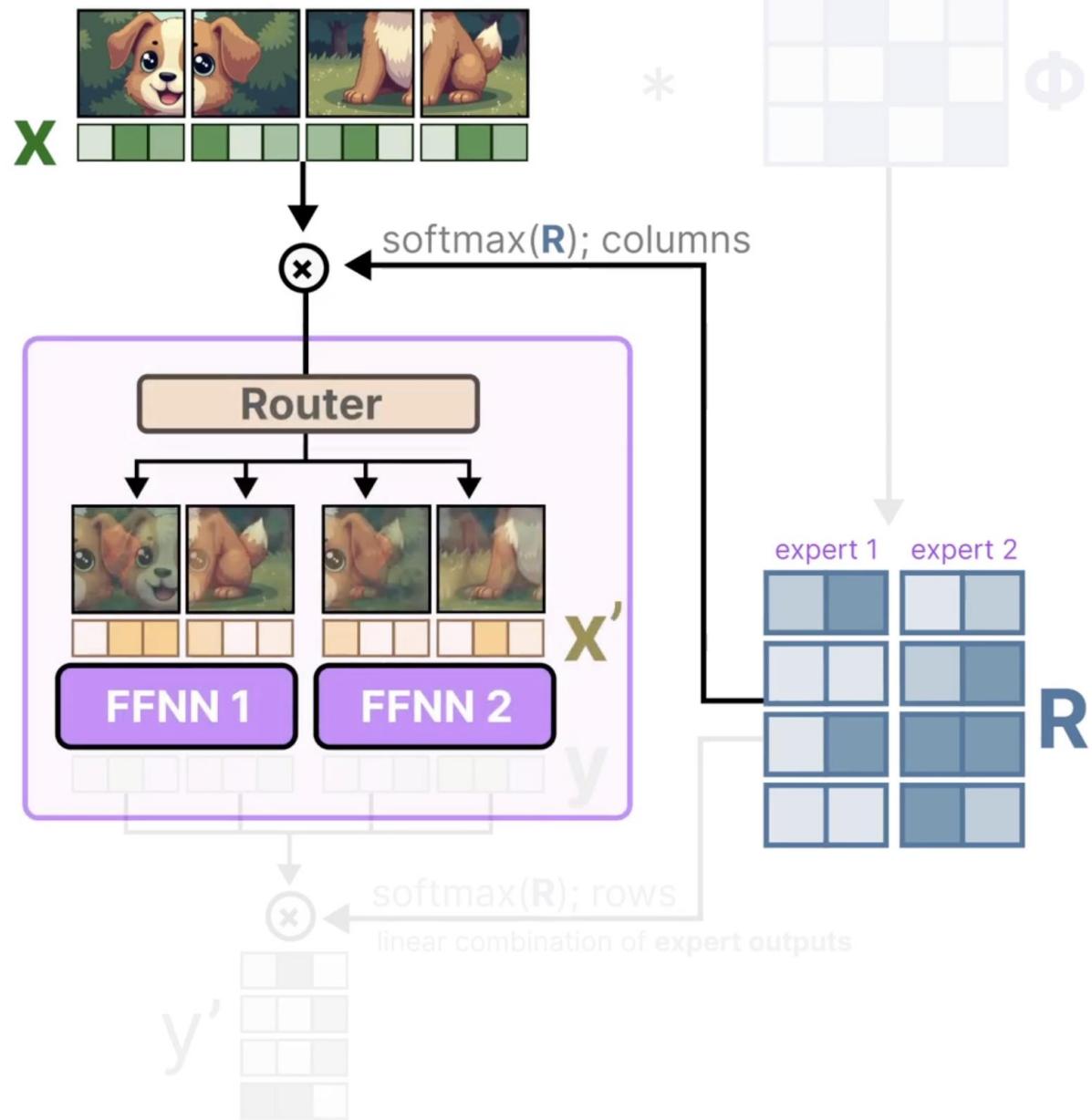
## From Sparse to Soft MoE



These are routed to the best suited **expert** for a particular patch input.



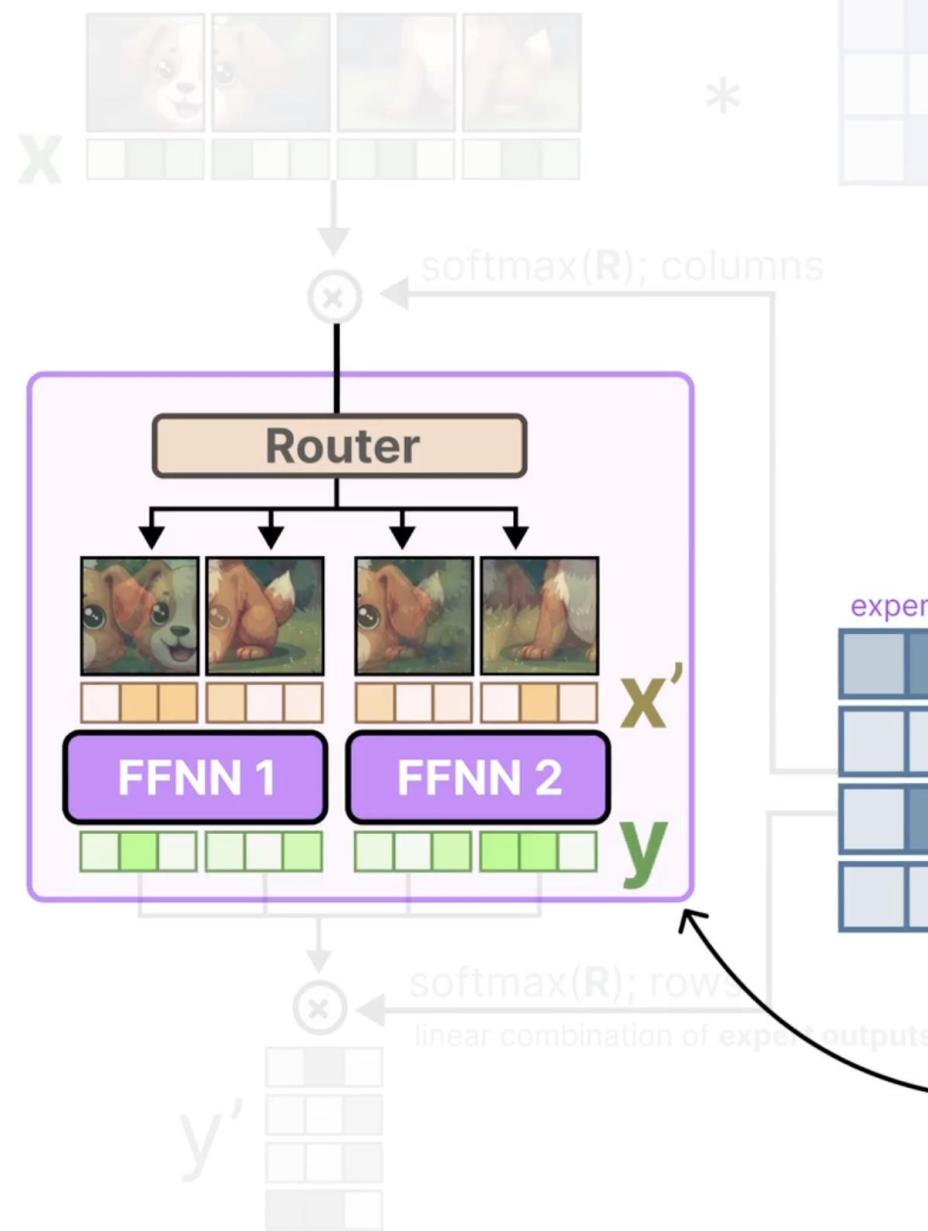
## From Sparse to Soft MoE



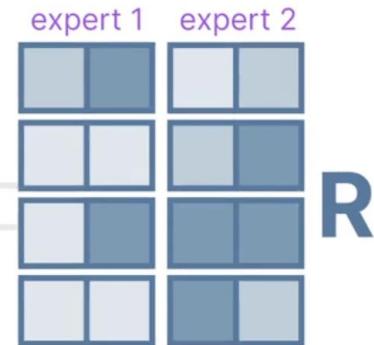
These are routed to the best suited **expert** for a particular patch input.

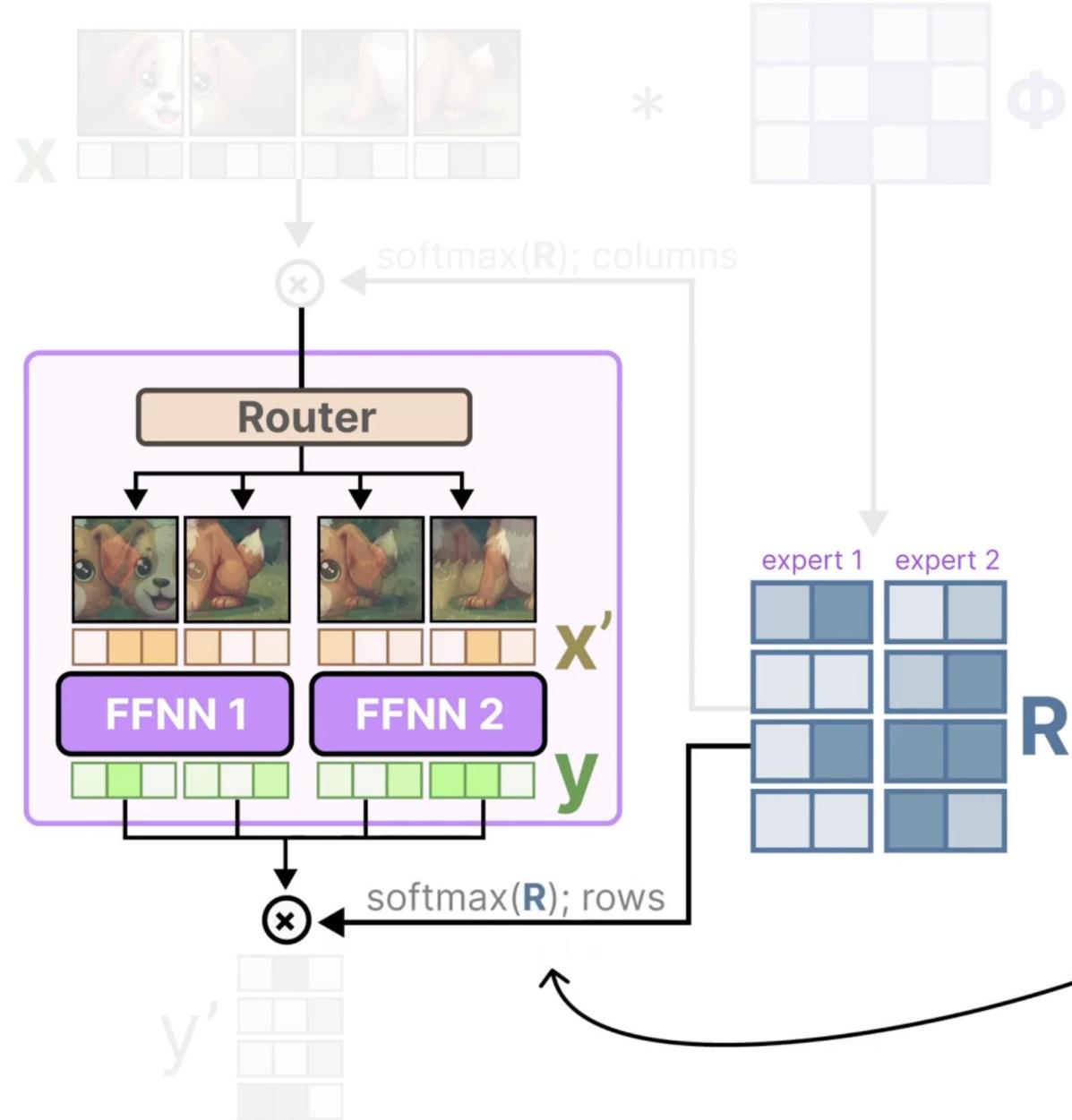


## From Sparse to Soft MoE



Each **expert** processes these linear combinations of the patch embeddings to generate  $y$ .

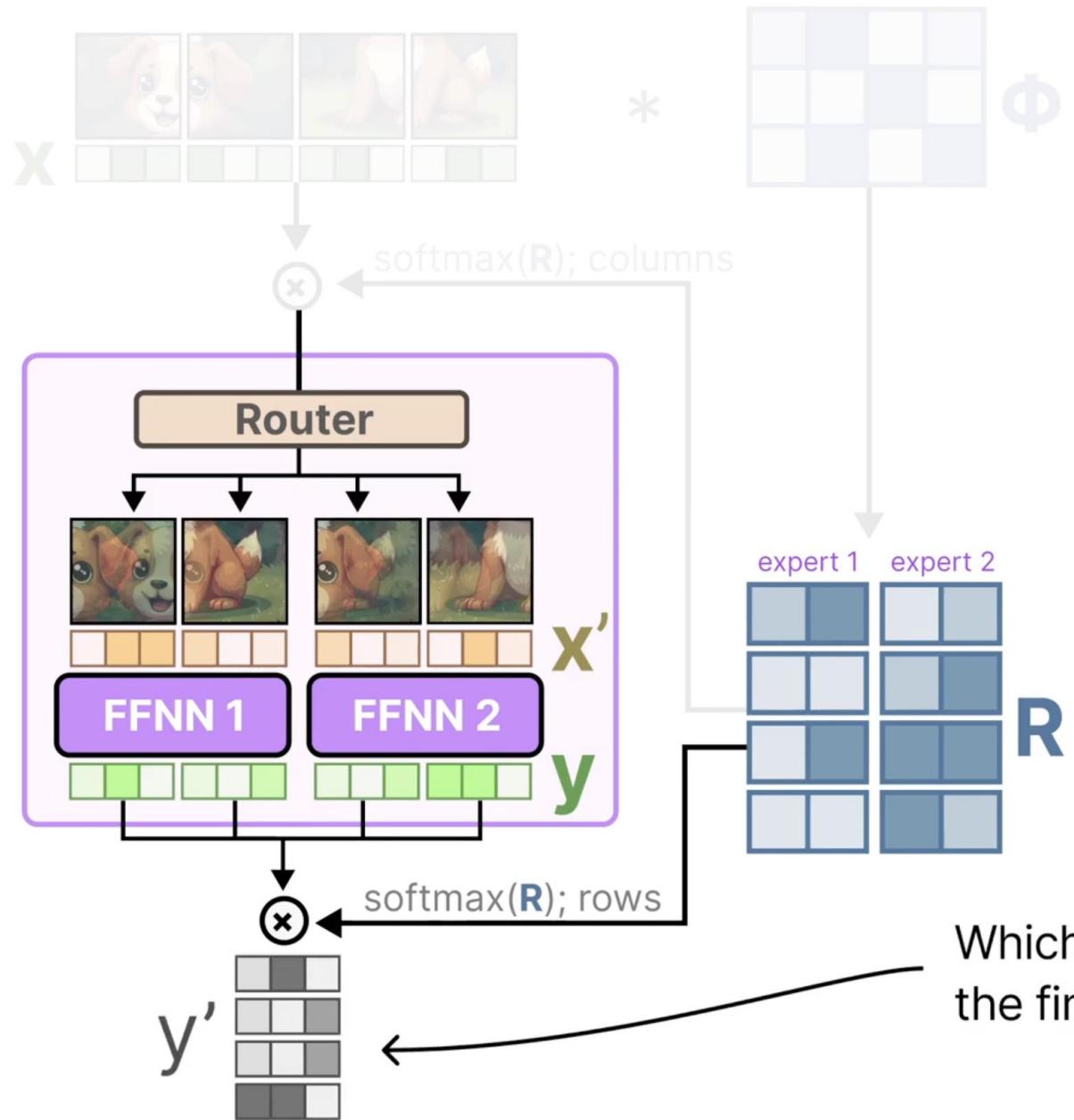




The output  $y$  is multiplied with the softmax of  $R$  but this time row-wise, creating a **linear combination of expert outputs** instead.



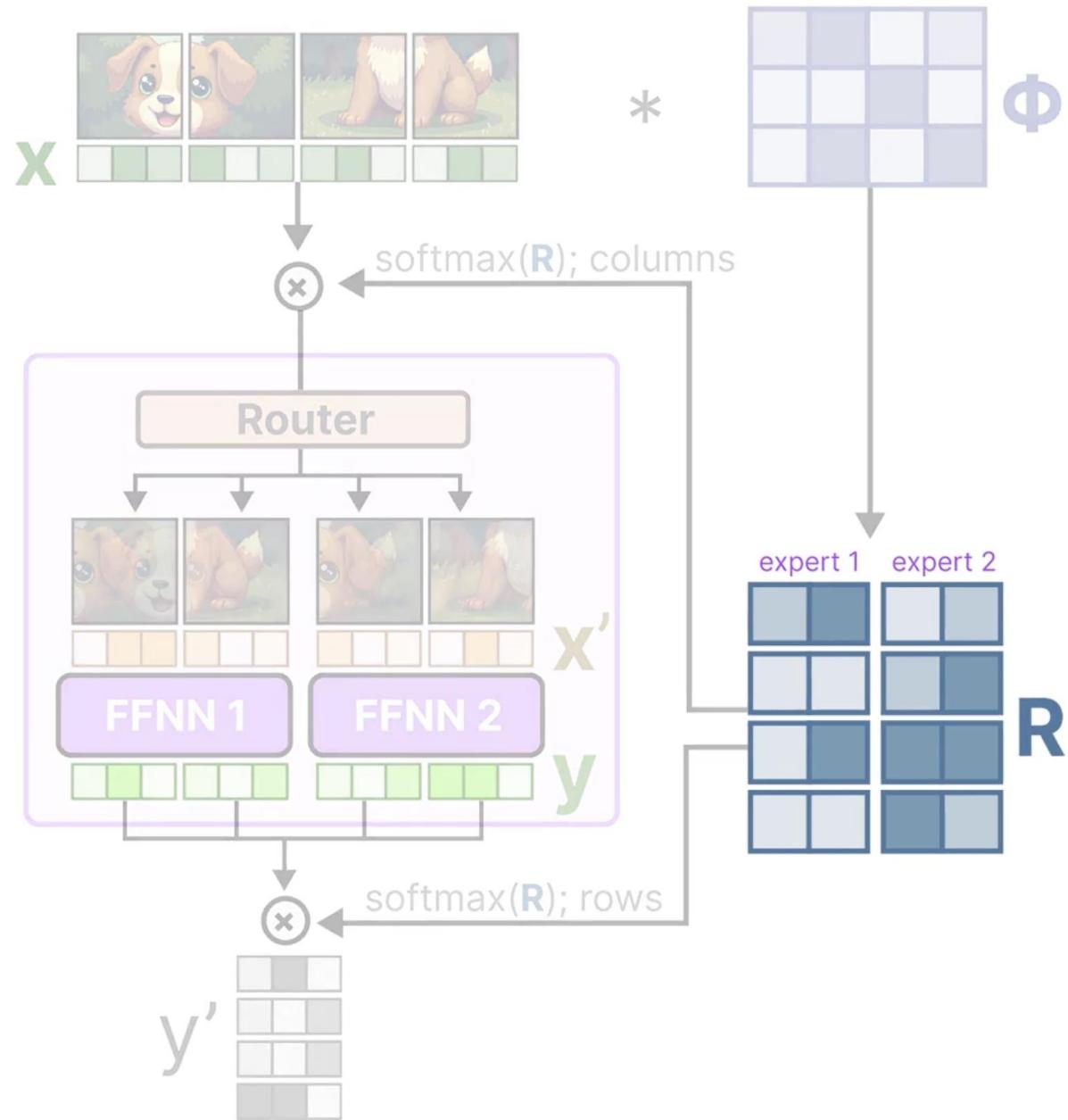
## From Sparse to Soft MoE



Which gives us  
the final output  $y'$ .



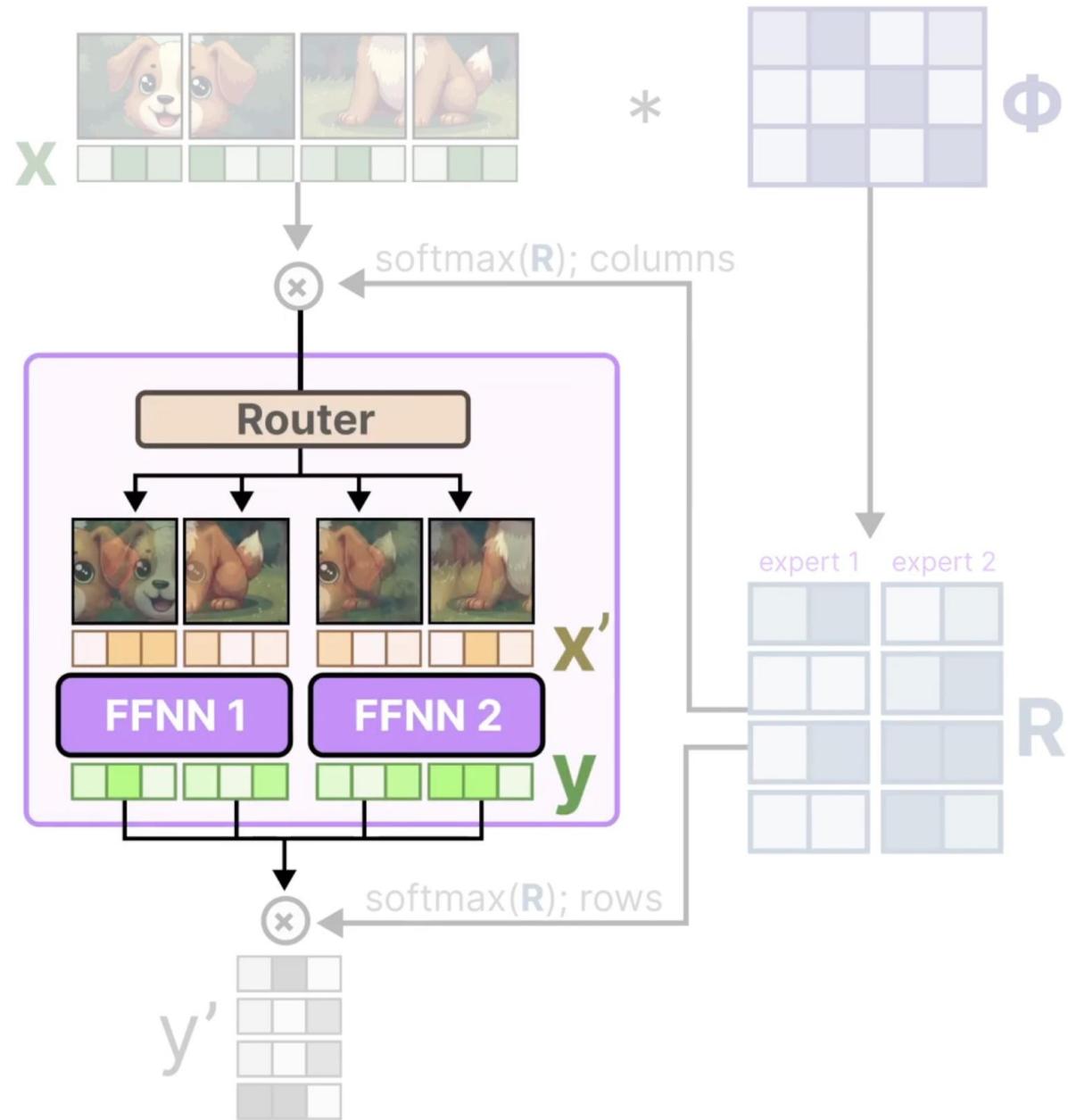
## From Sparse to Soft MoE



This architecture, together with the **router matrix**, affects the input on a **patch** level and the output on an **expert** level.



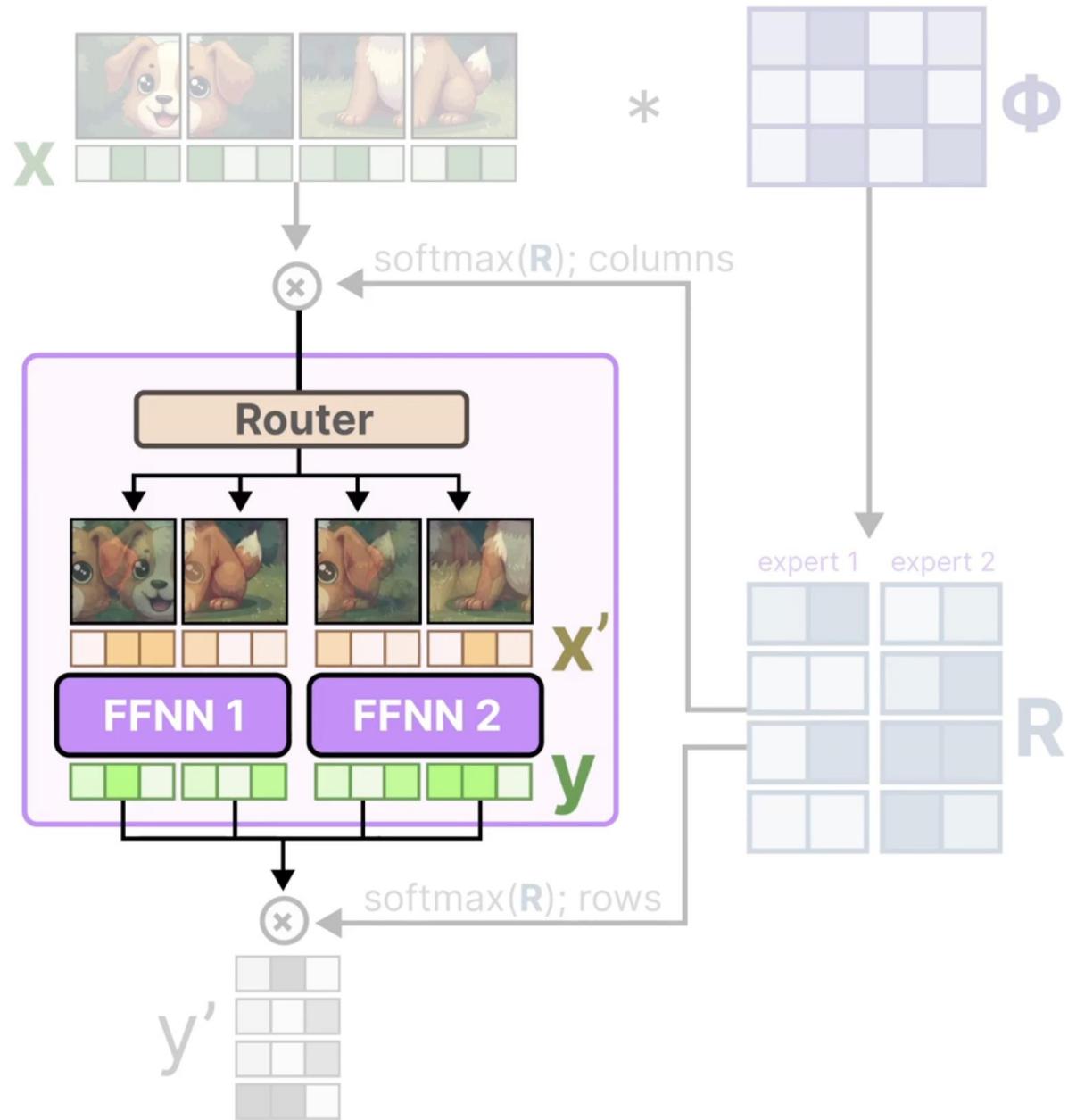
## From Sparse to Soft MoE



Since Transformers has found its way into the domain of vision, techniques like **MoE** are surprisingly transferable across domains.



## From Sparse to Soft MoE



So, keep in mind that whatever you learn about LLMs might also have theoretical and practical implications to **multimodal models!**





# Thank you

把AI系统带入每个开发者、每个家庭、  
每个组织，构建万物互联的智能世界

Bring AI System to every person, home and  
organization for a fully connected,  
intelligent world.

Copyright © 2024 XXX Technologies Co., Ltd.  
All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. XXX may change the information at any time without notice.



**ZOMI**

GitHub <https://github.com/chenzomi12/Allinfra>

# 引用与参考

- <https://www.youtube.com/watch?v=sOPDGQjFcuM&t=1s>
- <https://arxiv.org/abs/2308.00951>
- <https://arxiv.org/abs/2106.05974>
- PPT 开源在: <https://github.com/chenzomi12/AllInfra>

