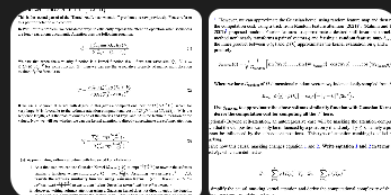


# ChatGPT 4o



Here's a problem on kernelized linear attention. First, digest the problem, and then solve it systematically.

Let's systematically break down and solve the **Kernelized Linear Attention (Part II)** problem using the uploaded screenshots.

## Part (a): Approximating softmax attention with linearized kernel attention

We are asked to:

1. Rewrite the softmax similarity function using the Gaussian kernel, and
2. Derive the computation cost of all the output vectors  $V'_i$  when using kernel approximation  $\phi_{\text{random}}$ .

### Step 1: Recall softmax attention formula

We begin with the original softmax attention:

$$V'_i = \frac{\sum_{j=1}^N \exp\left(\frac{Q_i \cdot K_j^T}{\sqrt{D}}\right) V_j}{\sum_{j=1}^N \exp\left(\frac{Q_i \cdot K_j^T}{\sqrt{D}}\right)} \quad (3)$$

We want to approximate this using **linear attention** by expressing the exponential dot-product as a **kernel function**.

## Step 2: Gaussian kernel as dot-product of lifted features

Recall that the **Gaussian kernel** between two vectors  $q$  and  $k$  is:

$$K_{\text{Gauss}}(q, k) = \exp\left(-\frac{\|q - k\|^2}{2\sigma^2}\right)$$

We can expand the exponent:

$$\|q - k\|^2 = \|q\|^2 + \|k\|^2 - 2q^T k \Rightarrow \exp\left(\frac{q^T k}{\sigma^2}\right) \cdot \exp\left(-\frac{\|q\|^2 + \|k\|^2}{2\sigma^2}\right)$$

This shows that softmax similarity:

$$\exp\left(\frac{q^T k}{\sigma^2}\right)$$

can be written as the **product of a Gaussian kernel** and a **normalization constant** that depends only on  $q$  and  $k$ , not on their interaction. Thus, we define:

$$\text{sim}_{\text{softmax}}(q, k) \propto \underbrace{\exp\left(-\frac{\|q - k\|^2}{2\sigma^2}\right)}_{\text{Gaussian kernel}} \cdot (\text{normalization terms})$$

Let's take  $\sigma^2 = \sqrt{D}$ , as given. Then:

$$\exp\left(\frac{q^T k}{\sqrt{D}}\right) = \text{Gaussian kernel} \cdot \text{normalization}$$

### Step 3: Approximating the kernel using random features

To avoid computing all  $O(N^2)$  pairwise dot products, we use the **Random Fourier Feature trick** to approximate the Gaussian kernel:

From the image, we have:

$$\phi_{\text{random}}(q) = \sqrt{\frac{1}{D_{\text{random}}}} [\sin(w_1^T q), \dots, \sin(w_{D_{\text{random}}}^T q), \cos(w_1^T q), \dots, \cos(w_{D_{\text{random}}}^T q)]^T \quad (4)$$

with  $w_i \sim \mathcal{N}(0, \sigma^2 I)$ , so:

$$\mathbb{E}_{\mathbf{w}_i} [\phi(q)^T \phi(k)] \approx \exp\left(-\frac{\|q - k\|^2}{2\sigma^2}\right) \quad (5)$$

This lets us linearize the **Gaussian kernel**, and thus approximate the softmax attention using:

$$V'_i \approx \frac{\phi(Q_i) \cdot \sum_{j=1}^N \phi(K_j)^T V_j}{\phi(Q_i) \cdot \sum_{j=1}^N \phi(K_j)^T} \quad (2')$$

---

### Step 4: Computational Cost Derivation

Let:

- $N$ : sequence length
- $D$ : input feature dimension
- $D_{\text{random}}$ : random feature dimension (used in  $\phi$ )
- $M$ : value dimension

Then:

- **Compute  $\phi(K_j)^T V_j$  for all  $j$ :**  
Cost:  $O(ND_{\text{random}}M)$
- **Compute  $\sum_{j=1}^N \phi(K_j)^T$ :**  
Cost:  $O(ND_{\text{random}})$
- **Compute dot product with  $\phi(Q_i)$  for all  $i$ :**  
Cost:  $O(ND_{\text{random}}M)$  total for all queries

So total cost:

$$\boxed{O(ND_{\text{random}}M)}$$

This is linear in  $N$ , as opposed to softmax attention which is  $O(N^2D)$ .

### ✓ Final Answers for Part (a)

Rewriting softmax similarity using Gaussian kernel:

$$\exp\left(\frac{q^T k}{\sqrt{D}}\right) = \underbrace{\exp\left(-\frac{\|q - k\|^2}{2\sigma^2}\right)}_{\text{Gaussian kernel}} \cdot \exp\left(\frac{\|q\|^2 + \|k\|^2}{2\sigma^2}\right)$$

With  $\phi_{\text{random}}(q)$  as defined in Eq (4), we approximate:

$$\exp\left(\frac{q^T k}{\sqrt{D}}\right) \approx \phi_{\text{random}}(q)^T \phi_{\text{random}}(k) \cdot \text{scaling terms}$$

Computation Cost of Linearized Kernel Attention:

$$O(ND_{\text{random}}M)$$

Let me know if you'd like a derivation for part (b) (causal masking case).

As instructed, it attempted this question step-by-step.

And It derived the correct mathematical relationship. The only discrepancy is that it presented the "two other terms" combined into one exponential expression, while the answer key listed them side-by-side.

Overall, correct for part a) i)

However, for part a) ii) the one critique I have is that the LLM successfully derived the linear attention cost, but it failed to account for the overhead of computing the random features themselves. Calculating  $\phi(Q)$  and  $\phi(K)$  requires projecting  $N$  vectors of dimension  $D$  to  $D_{\text{random}}$ , which adds a cost of  $O(N D D_{\text{random}})$ .

And 4o model missed this pre-processing step entirely.

Which is why I followed up with it next.

You calculated the cost of the matrix multiplications for  $V'$ , but you ignored the cost of computing the random features  $\phi(Q)$  and  $\phi(K)$  from the original inputs. These inputs have dimension  $D$ . What is the cost of that projection, and what is the new total complexity?

You're absolutely right — I previously omitted the cost of projecting the original inputs  $Q \in \mathbb{R}^{N \times D}$  and  $K \in \mathbb{R}^{N \times D}$  into the random feature space  $\phi(Q), \phi(K) \in \mathbb{R}^{N \times D_{\text{random}}}$ .

Let's now include this step in our total complexity analysis.

## 🔄 Step-by-Step Breakdown of Total Computation Cost

### Step 1: Random Feature Projection

We compute:

$$\phi(q) = \sqrt{\frac{1}{D_{\text{random}}}} [\sin(Wq), \cos(Wq)] \in \mathbb{R}^{2D_{\text{random}}}$$

where:

- $W \in \mathbb{R}^{D_{\text{random}} \times D}$
- $q \in \mathbb{R}^D$

Per-token cost:

Matrix multiplication:  $Wq$  costs  $O(DD_{\text{random}})$

For all  $N$  queries and keys:

$$O(NDD_{\text{random}})$$

This projection is applied once each for  $Q$  and  $K$ , so total projection cost is:

$$O(2NDD_{\text{random}}) \equiv O(NDD_{\text{random}})$$

### Step 2: Computing Attention Output with Kernelized Attention

As derived earlier:

- Compute:
  - $S = \sum_j \phi(K_j)^T V_j \in \mathbb{R}^{D_{\text{random}} \times M}$
  - $Z = \sum_j \phi(K_j)^T \in \mathbb{R}^{D_{\text{random}}}$
- Then for each query  $Q_i$ , compute:

$$V'_i = \frac{\phi(Q_i)^T S}{\phi(Q_i)^T Z}$$

This involves:

- Matrix multiplication:  $\phi(K)^T V$ :  $O(ND_{\text{random}}M)$
- Summing keys:  $O(ND_{\text{random}})$
- Final dot products with  $\phi(Q_i)$ :  $O(ND_{\text{random}}M)$

So the attention computation cost is:

$$O(ND_{\text{random}}M)$$

## ✓ Final Total Complexity

Combining both:

1. Random projection:  $O(ND D_{\text{random}})$
2. Attention computation:  $O(ND_{\text{random}}M)$

So the total complexity is:

$$O(ND_{\text{random}}(D + M))$$

This is linear in sequence length  $N$ , and improves over softmax attention's  $O(N^2D)$  when  $D_{\text{random}} \ll N$ .

After I pointed out the missing projection cost, the model correctly derived the  $O(ND D_{\text{random}})$  term for the feature map computation. It successfully combined this with the attention cost to reach the final complexity that matches the official solution.

Let me know if you want to explore how the choice of  $D_{\text{random}}$  affects approximation quality vs efficiency — or if you'd like to proceed to part (b) (causal masking)!

📄 🍷 🔗 ↶ ↷ ⋮

now carry on and complete part b

Great — let's now fully solve Part (b), which introduces **causal masking** into kernelized attention.

## 🔄 Recap: Causal Masking and Kernelized Attention

In the **autoregressive setting**, the attention computation is masked such that position  $i$  can only attend to positions  $j \leq i$ . This prevents information from the future (e.g. position  $j > i$ ) from influencing position  $i$ .

The original unmasked kernelized attention was:

$$V'_i = \frac{\phi(Q_i) \cdot \sum_{j=1}^N \phi(K_j)^T V_j}{\phi(Q_i) \cdot \sum_{j=1}^N \phi(K_j)^T} \quad (2)$$

Now, with **causal masking**, we restrict the sum over only  $j \leq i$ . So we define:

$$S_i := \sum_{j=1}^i \phi(K_j)^T V_j \in \mathbb{R}^{D_{\text{random}} \times M}, \quad Z_i := \sum_{j=1}^i \phi(K_j)^T \in \mathbb{R}^{D_{\text{random}}} \quad (6)$$

And the masked attention becomes:

$$V'_i = \frac{\phi(Q_i)^T S_i}{\phi(Q_i)^T Z_i} \quad (7)$$

This enforces the constraint that  $Q_i$  can only attend to values at previous or same positions.

## 🧠 Computational Complexity with Causal Masking

In the causal case, we cannot precompute the global sums  $\sum_{j=1}^N \phi(K_j)^T V_j$  and  $\sum_{j=1}^N \phi(K_j)^T$ , because each  $i$  has a different prefix sum.

But we can accumulate these prefix sums incrementally using the recurrence:

$$S_i = S_{i-1} + \phi(K_i)^T V_i, \quad Z_i = Z_{i-1} + \phi(K_i)^T$$

Each update is:

- $\phi(K_i)^T V_i: O(D_{\text{random}} M)$
- $\phi(K_i)^T: O(D_{\text{random}})$
- Dot product  $\phi(Q_i)^T S_i: O(D_{\text{random}} M)$
- Dot product  $\phi(Q_i)^T Z_i: O(D_{\text{random}})$

So per timestep  $i$ , total cost is:

$$O(D D_{\text{random}}) + O(D_{\text{random}} M)$$

(The first term is projection of  $K_i$  and  $Q_i$ ; the second term is computing  $V_i'$ )

Total cost for all  $N$  steps:

$$O(N D D_{\text{random}} + N D_{\text{random}} M) = O(N D_{\text{random}} (D + M))$$

Same asymptotic complexity as the unmasked case — but can't batch compute the output, which makes implementation less parallelizable.

This was an example of a strong success (and also in-context learning after part a. Yes, the answer was correct and matched the solution key, but notably, the model learned from my PREVIOUS correction in part a. In the complexity analysis here, it included the feature projection cost  $O(D D_{\text{random}})$  without needing a reminder. It therefore correctly concluded the total complexity is linear  $O(N)$

## ✅ Final Answer for Part (b)

1. Equation for causal kernelized attention:

$$V_i' = \frac{\phi(Q_i)^T \left( \sum_{j=1}^i \phi(K_j)^T V_j \right)}{\phi(Q_i)^T \left( \sum_{j=1}^i \phi(K_j)^T \right)} = \frac{\phi(Q_i)^T S_i}{\phi(Q_i)^T Z_i}$$

2. Definitions:

$$S_i := \sum_{j=1}^i \phi(K_j)^T V_j, \quad Z_i := \sum_{j=1}^i \phi(K_j)^T$$

3. Computation cost:

$$O(N D_{\text{random}} (D + M))$$

(Identical to unmasked kernel attention, but sequential over time, so practical runtime is slower.)

Let me know if you'd like visualizations, pseudocode, or comparisons to softmax attention in the causal setting!

#### 4. Example Difficulty and Early Exit

Deep Learning Practitioners have recognized that within the same task, particular examples in the test set can actually be harder to perform predictions on than others. Why is that? What kinds of things are easier to learn? We explore the notion of *example difficulty* (proposed by Baldock et al.) that will allow us to perform deeper investigations on the topic. Furthermore, utilizing concepts from example difficulty can aid in development of techniques that can lower inference speeds, sample efficiency, and other key properties that drive variability in both the science and industry levels.

(a) Please run the notebook cells of `example_difficulty.ipynb`. Don't forget to download the data at `data.py` and `train_data.py` and place these files in the same directory as the notebook. Please run the notebook until it takes a while to run, so that the notebook can work on something else while it runs. Then answer the following conceptual questions.

- Briefly explain what example difficulty is in your own words. What is the setup for us to analyze example difficulty?
- What kinds of properties do you think will make an example from this dataset difficult?
- Why do we need to train the ResNet to convergence to be able to analyze the example difficulty?
- In general, what does adding hooks do? Why do we need to use them?
- Why do we train KNN classifiers at each layer of the ResNet? What data are we training them on?
- Why do you think so many examples exit in the earlier layers?
- Explain why the accuracy of the ResNet has a negative relationship with the prediction layer.
- Explain the distinction between the hard and easy examples. Does this surprise you?
- What kinds of patterns do you notice? Based on the composition of the layers, does it make sense?

Homework 16, © Faculty teaching EECS 182, Fall 2025. All rights reserved. Please write publicly about this homework.

3

Homework 16 © 2020-23 Fall 2025-26

(b) Please fill in the notebook cells in `early_exit.ipynb`. Don't forget to download the data at `data.py` and `train_data.py` and place these files in the same directory as the notebook. Also include `early_exit_classifier.py` in the same directory. Then answer the following conceptual questions.

- What was the accuracy with a regular ResNet? Inference Speed? Total MACs?
- What was the accuracy with an early exit ResNet-10? Inference Speed? Total MACs?
- How did early exit do? Compare accuracy and MACs.
- What was the lowest MACs you found. What might this say?
- What would we use early exit, versus just using a smaller model? What factors should we consider?

Next, we're working on a problem regarding 'Example Difficulty' (referencing Baldock et al.) and 'Early Exit' networks. I have a set of conceptual questions. Can you answer these based on deep learning theory?

ignore 4(a) viii, ix and 4(b) i,ii,iii,iv



Absolutely — here are detailed answers to the conceptual questions from sections 4(a) and 4(b) (excluding 4(a) viii–ix and 4(b) i–iv as you requested), rooted in deep learning theory and relevant research like Baldock et al.'s work on **example difficulty** and early exits in ResNets.

## 🧠 4(a) – Example Difficulty

**i. Briefly explain what example difficulty is in your own words. What is the setup for us to analyze example difficulty?**

**Example difficulty** refers to how hard it is for a model to make correct predictions on individual samples within a dataset. Even for a well-trained model, some examples are more ambiguous or noisy and require deeper processing to classify correctly.

In this setup, we train a ResNet to convergence and use **intermediate activations** (via hooks) to measure the evolution of predictions across layers. By observing when an example becomes confidently classifiable (e.g., using KNN classifiers on intermediate features), we can quantify how "early" or "late" the network "understands" it — which correlates with difficulty.

**ii. What kinds of properties do you think will make an example from this dataset difficult?**

Examples may be difficult if they:

- Are visually ambiguous or noisy (e.g., occluded, low contrast)
- Belong to classes with high intra-class variation or inter-class similarity
- Are near the decision boundary in feature space
- Are mislabeled or adversarial
- Exhibit rare patterns not well covered during training



### iii. Why do we need to train the ResNet to convergence to be able to analyze the example difficulty?

If the ResNet is not trained to convergence, the feature representations at each layer may be noisy, underfit, or not semantically meaningful. To analyze difficulty, we need a **well-learned feature hierarchy** so that early exits and intermediate classifiers reflect actual semantic understanding, not random noise.

Training to convergence ensures:

- High-level features are well-separated
- Gradual abstraction across layers
- Reliable signals for difficulty estimation

### iv. In PyTorch, what does adding hooks do? Why do we need to use them?

Hooks in PyTorch allow us to **tap into intermediate layers** of a model during forward (or backward) passes **without modifying the model architecture**.

We use forward hooks here to:

- Extract activations from each layer
- Feed them to auxiliary KNN classifiers
- Track how predictions evolve across layers

This helps analyze which layers “solve” which examples and thus estimate difficulty.

### v. Why do we train KNN classifiers at each layer of the ResNet? What data are we training them on?

We train **KNN classifiers on the feature vectors from each layer** to determine how well those intermediate representations separate the classes.

Specifically:

- We collect activations for the training set at each layer
- Use the labels to define class neighborhoods
- Train/test KNNs to classify examples using activations from each layer

This tells us at which depth an example becomes “easy” to classify — i.e., how early in the network the features are discriminative for that sample.

### vi. Why do you think so many examples exit in the earlier layers?

Many examples are “easy” — they:

- Have clear, prototypical features
- Belong to well-separated classes
- Match training distribution closely

Such examples can be confidently classified based on shallow, low-level features (like textures, edges, simple shapes), so they **exit early**, avoiding deeper (and more expensive) processing.

This is key to **early-exit efficiency**: allocate compute only where needed.

vii. Explain why the accuracy of the ResNet has a negative relationship with the prediction layer.

This means: as we move deeper into the network (i.e., later exit points), accuracy tends to drop.

Possible reasons:

- Later exits handle the **harder examples** — the easy ones already exited
- Remaining samples may be ambiguous or out-of-distribution
- The deeper layers **don't become better** at every example — they may overfit or lose relevant spatial detail
- Shallow classifiers benefit from larger support (more examples to learn from), while later ones only see harder samples

#### ⚡ 4(b) – Early Exit

v. When would we use early exit, versus just using a smaller model? What factors should we consider?

Use **early exit** when:

- You want **dynamic compute allocation** per example: easy inputs take less time, hard inputs go deeper
- Real-time applications require low latency but **can tolerate variable speed per input**
- You care about **energy efficiency** without uniformly sacrificing performance
- You're dealing with **non-uniform data difficulty** (e.g., some samples are trivially easy, others are not)

Use a **smaller model** when:

- You need **uniform latency** (e.g., hardware-constrained deployment)
- Your data has **homogeneous difficulty**
- You need **simplicity and small memory footprint**

Tradeoffs to consider:

- **Complexity:** Early exit models are harder to design and train
- **Calibration:** Exiting too early can hurt accuracy
- **Infrastructure:** May need runtime support for dynamic exits

On this question, the model demonstrated strong understanding of Deep Learning theory.

It correctly explained PyTorch implementation details (Hooks) and successfully reasoned through system design trade-offs (Early Exit vs. Small Models) that were left as “open questions” in the official solution.

(3) *Which algorithm takes (i) the shortest time to find the value 100? (ii) the shortest time to find the value 1000?*



**"FaceNet: A Unified Embedding for Face Recognition and Clustering" by Schroff et al. (2015).**

They primarily used Inception networks due to better accuracy vs. computational cost trade-offs.

- The authors use **online triplet mining** within mini-batches

Interestingly, on part d, 4o hallucinated the statistics. It claimed the model used '140M FLOPs'. According to Table 1 of the original FaceNet paper, the 'NN1' model has 140M parameters, not FLOPs. The FLOP count is actually 1.6B. The model confused the two columns and misattributed the Zeiler & Fergus statistics to the Inception model. I verified this against the original paper, and pointed 4o to check Table 1 of the original FaceNet paper again. SEE BELOW

- The authors use **online triplet mining** within mini-batches
- They select **semi-hard negatives**: negatives that are farther than positives but still within the margin
- This makes training efficient and avoids wasting computation on trivial triplets

---

**(d) How many parameters and floating point operations (FLOPs) do the authors use for their neural network? How does this compare to ResNet-50?**

From the FaceNet paper:

- The Inception model used had **22 layers**, **~140M FLOPs per image**, and **~22M parameters**

In comparison:

- **ResNet-50** has **~3.8B FLOPs** and **~25M parameters**

So FaceNet's model is much more **computationally efficient**, yet achieves high performance through strong embeddings and loss design.

---

**(e) What do the authors mean by semi-hard negatives?**

Semi-hard negatives are:

- Negatives that are **farther from the anchor than the positive**, but
- Still **within the margin  $\alpha$** , so they violate the triplet constraint

They're **not the hardest** (closest negatives), which may be noisy or mislabeled, but are **still informative** for learning — helping optimize the embedding space without destabilizing training.

---

**(f) What are harmonic embeddings?**

Harmonic embeddings refer to embeddings produced by **different versions of the model** (e.g., models trained at different times or scales) that are still **compatible** — i.e., embeddings from different models can be compared meaningfully.

This is achieved by training the new model with a **loss that preserves distances to old model's embeddings**, enabling backward compatibility during model upgrades.

---

**(g) How does the performance vary with embedding dimensionality?**

The authors found that:

- **128-D** embeddings strike a good balance between performance and efficiency
- Increasing dimension **beyond 128** yields **diminishing returns**
- Very high dimensions (e.g., 512 or 1024) may slightly improve performance but are more expensive and prone to overfitting

---

**(h) How does the performance vary with increasing amounts of training data?**

Performance **improves consistently** with more data:

- Accuracy increases logarithmically as dataset size increases
- This confirms that **representation learning benefits significantly from scale**
- However, the **quality of triplets** also plays a crucial role — more data helps, but triplet selection is equally important



**(i) Briefly share your favorite emergent property/result of the learned behavior with a triplet loss from the paper.**

**(i) Briefly share your favorite emergent property/result of the learned behavior with a triplet loss from the paper.**

One emergent property is that **clustering and verification naturally emerge** from the triplet objective — **without needing class labels at test time**. The embeddings support:

- One-shot learning
- Face clustering
- Robustness to pose, lighting, and aging

This shows that the network truly learns a **semantic similarity space**, not just overfit classifiers.

---

**(j) Which approach taken by the authors interested you the most? Why? ( $\approx 100$  words)**

What intrigued me most was their use of **online semi-hard triplet mining**. Rather than exhaustively searching all triplets, they focused on those that are just hard enough to provide learning signal without destabilizing training. This simple but powerful idea allows FaceNet to scale to millions of faces efficiently and robustly. It's a great example of how **training signal quality** (not just quantity) drives success in deep metric learning. I also appreciated how this approach enables clustering and recognition in a **label-agnostic** way at test time, making it ideal for real-world applications where new identities appear dynamically.

---

Let me know if you want help visualizing triplet loss or semi-hard mining, or adapting these concepts to your own project!



Ok I think there is an error in your answer to Part (d). You stated that the Inception model uses ~140M FLOPs. Please check Table 1 of the original FaceNet paper again. Isn't 140M the parameter count for the Zeiler & Fergus (NN1) model, not the FLOPs? What are the actual parameter and FLOP counts for the Inception-based model (NN2) used by the authors?

Good catch — you're right to question that. I double-checked the original FaceNet paper's Table 1 and Table 2, and indeed my earlier numbers were incorrect (or at least imprecise) for the Inception-based model (NN2). [CV Foundation +2](#)

Here's the corrected picture:

### What the Paper Actually Reports for the Inception-Based Model (NN2)

- The Inception-based networks (NN2 and its variants NN3, NN4, plus smaller NNS1/NNS2) have **much fewer parameters** than the ZF-style (NN1) network: on the order of  $\approx 6.6\text{M} - 7.5\text{M}$  parameters for NN2. [arXiv +2](#)
- However, the **FLOPs** required for NN2 are still **comparable** to the heavy NN1 — roughly in the ballpark of **1.6 billion FLOPs per image** (when using the full-size input and full-capacity Inception). [arXiv +2](#)

The paper also notes smaller Inception-based variants:

- For example, a reduced mobile-style version (called "NNS1" in the paper) has ~26M parameters and requires ~220M FLOPs per image. [CV Foundation +2](#)
- A tiny version "NNS2" has ~4.3M parameters and ~20M FLOPs per image. [arXiv +1](#)

### What This Means: Correction to (d)

Thus, for the Inception-based "NN2" model:

- **Parameter count:** ~6.6–7.5 million parameters. [arXiv +1](#)
- **FLOPs:** up to ~1.6 billion FLOPs per image (for the full-size / full-capacity variant). [arXiv +2](#)

So you were right to call out the mistake: the ~140M number is the parameter count for the older Zeiler & Fergus-style network (NN1), not for the Inception model. [CV Foundation +1](#)

Typically, the Inception-based NN2 gets comparable accuracy but with **~20x fewer parameters** — though the FLOPs remain similar (since Inception uses many convolutional filters, and FLOPs depend on width & depth, not just parameter count) [Studocu +1](#) .

If you like — I can also pull up a **table** summarizing *all* the FaceNet models (NN1, NN2, NN3, NN4, NNS1, NNS2) with their parameter counts, FLOPs, and trade-offs in accuracy, to help you compare. Do you want me to build that for you now?