

MicroGPT

Visual Walkthrough

A complete GPT built from scalar autograd in ~200 lines of Python
No tensors. No PyTorch. No CUDA. Just math.

AI Engineering Cohort | TensorTonic
Teaching reference for 200 engineers

Model configuration

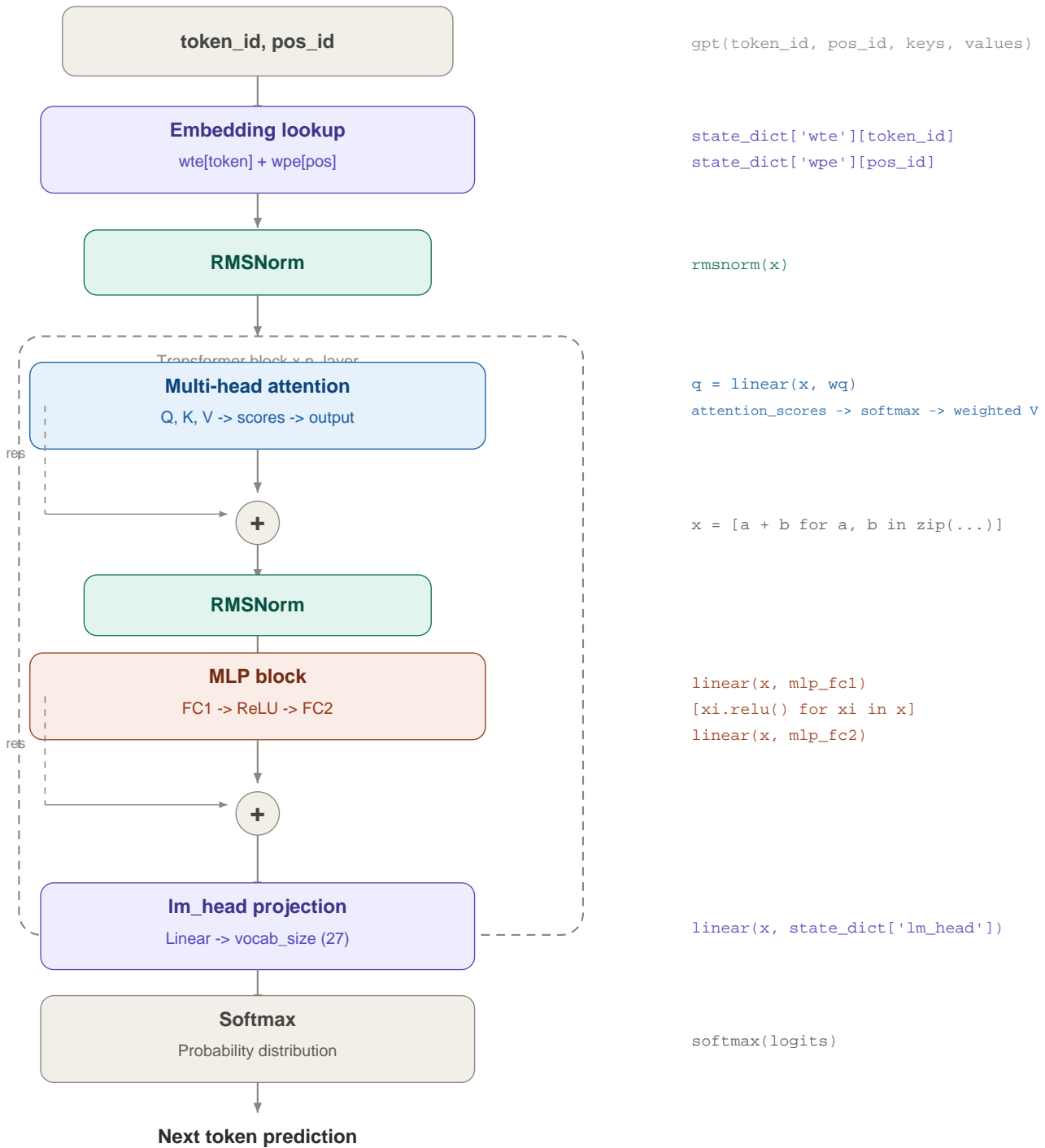
Concept	Value	Code
Vocabulary	27 chars (a-z + BOS)	<code>chars = sorted(set('').join(docs))</code>
Embedding dim	16	<code>n_embed = 16</code>
Attention heads	4 heads x 4-dim each	<code>n_head = 4, head_dim = n_embed // n_head</code>
Layers	1 transformer block	<code>n_layer = 1</code>
Context window	16 tokens max	<code>block_size = 16</code>
Total params	~5K Value objects	<code>params = [p for matrix in state_dict.values() ...]</code>
Batch size	8 names per step	<code>batch_size = 8</code>
Training steps	20	<code>num_steps = 20</code>
Optimizer	Adam (B1=0.85, B2=0.99)	<code>LR decays linearly 0.01 -> 0</code>
Init	Gaussian(0, 0.08)	<code>random.gauss(0, 0.08)</code>

What's in this document

1. End-to-end GPT architecture — the `gpt()` function as a pipeline
2. Autograd / Value class — scalar computation graph + `backward()`
3. Multi-head attention — Q, K, V projections, head splitting, KV cache
4. Training loop + Adam optimizer — forward, loss, backward, update
5. Embedding + positional encoding — token + position lookup tables

1. GPT Architecture Overview

The gpt() function — decoder-only transformer, one token at a time

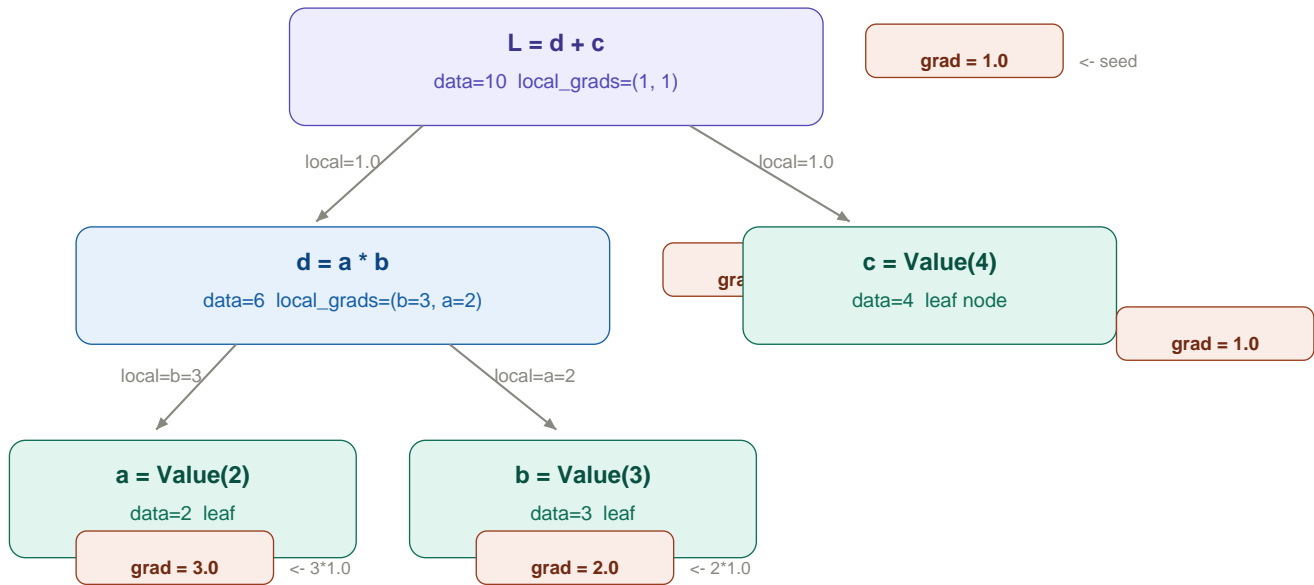


Key insight: gpt() processes one token at a time. The KV cache (keys, values lists) stores previous tokens' representations so attention can look backward without recomputing. This is how the model 'remembers'.

2. Autograd — The Value Class

Every operation builds a DAG node: `Value(data, children, local_grads)`

Example computation: $L = (a * b) + c$ where $a=2, b=3, c=4$



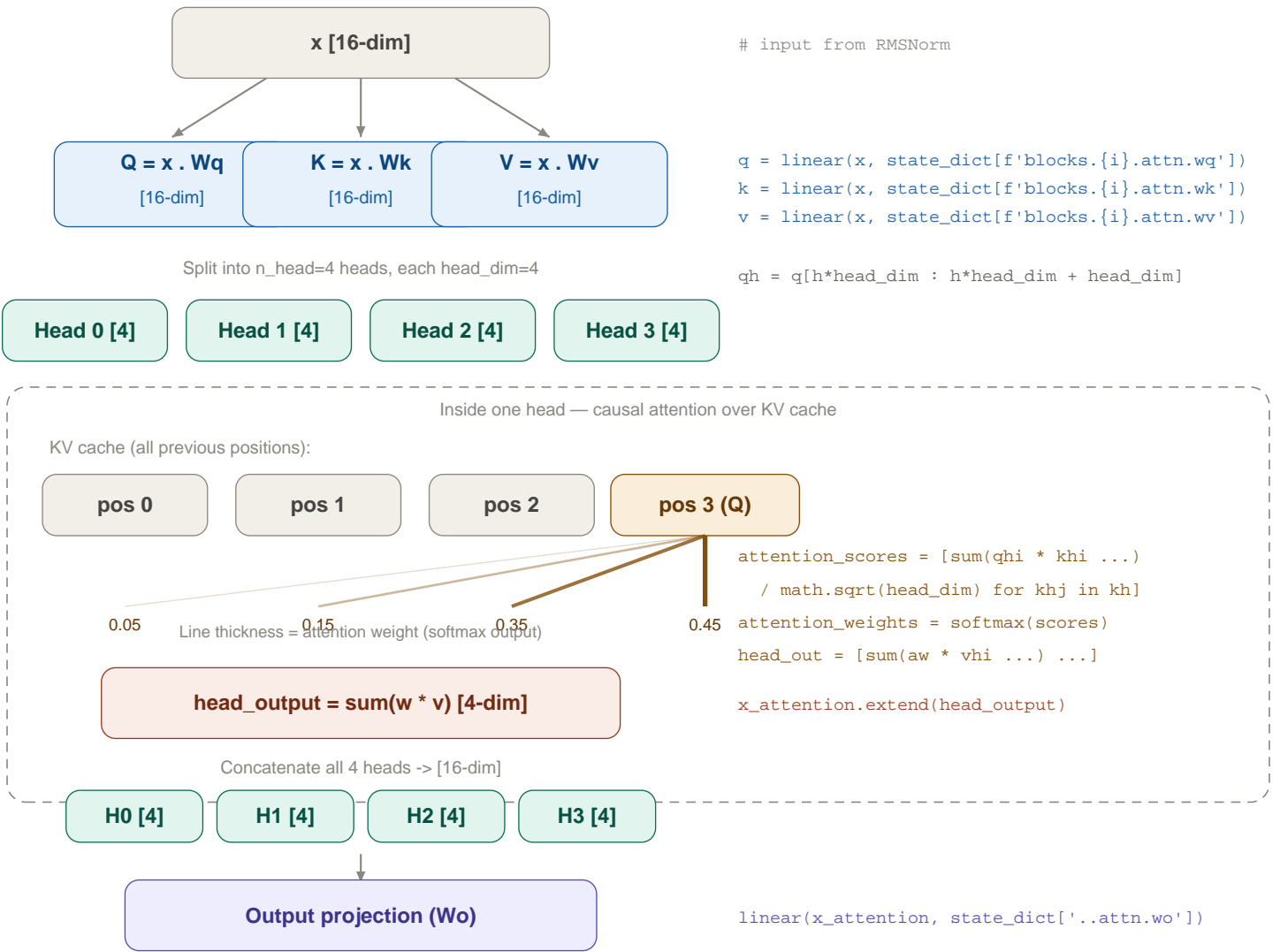
backward() — chain rule in 3 lines:

```
for v in reversed(topological_graph):
    for child, local_grad in zip(v._children, v.local_grads):
        child.grad += local_grad * v.grad
```

Operation	Forward	local_grads	Chain rule intuition
<code>a + b</code>	<code>a.data + b.data</code>	<code>(1.0, 1.0)</code>	Addition passes gradient through unchanged
<code>a * b</code>	<code>a.data * b.data</code>	<code>(b.data, a.data)</code>	Multiply swaps: $dL/da = b * \text{upstream}$
<code>a ** n</code>	<code>a.data ** n</code>	<code>(n * a^(n-1),)</code>	Power rule from calculus
<code>a.log()</code>	<code>log(a.data)</code>	<code>(1/a.data,)</code>	$d/da \log(a) = 1/a$
<code>a.exp()</code>	<code>exp(a.data)</code>	<code>(exp(a.data),)</code>	Exponential is its own derivative
<code>a.relu()</code>	<code>max(0, a.data)</code>	<code>(float(a > 0),)</code>	Gradient is 1 if positive, 0 otherwise

Key insight: This is the exact same algorithm PyTorch uses internally — topological sort + chain rule. The only difference is PyTorch operates on tensors (batched), while `microgpt.py` operates on individual floats (scalar).

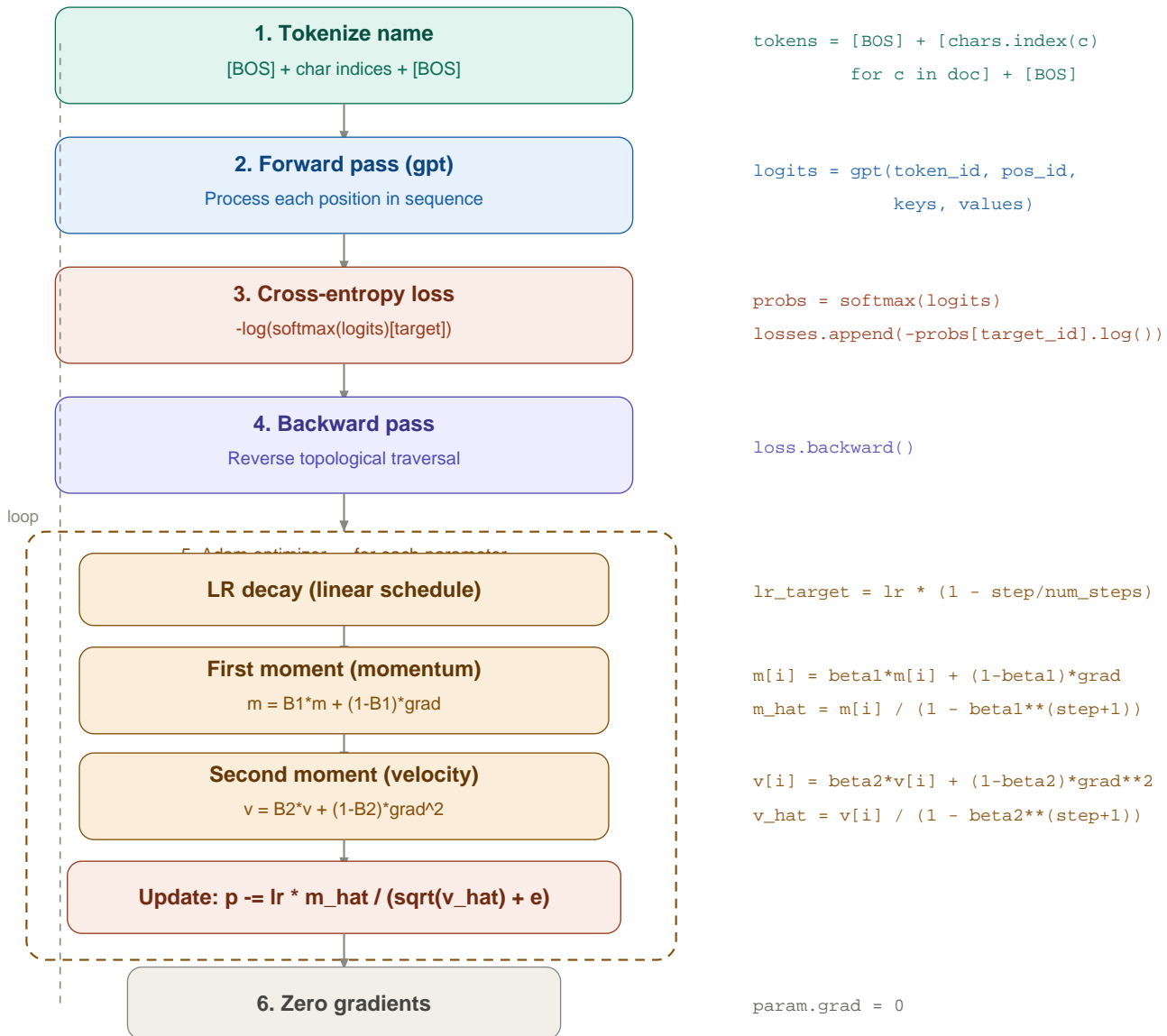
3. Multi-Head Attention



KV cache: `keys[i].append(k)` and `values[i].append(v)` store every token's K and V. At inference, the model doesn't recompute past tokens — it just extends the cache. This is exactly how production LLMs work.

4. Training Loop + Adam Optimizer

One training step (repeated 20 times, batch_size=8):



Why Adam? First moment (m) = exponential moving average of gradients (direction). Second moment (v) = EMA of squared gradients (magnitude). Bias correction compensates for initialization at zero in early steps.

Architecture flow

Input	token_id (0-26), pos_id (0-15)
Embed	wte[token] + wpe[pos] -> [16-dim]
Norm	RMSNorm(x)
Attn	Q,K,V -> 4 heads -> concat -> Wo
Residual	x = attn_out + x_residual
Norm	RMSNorm(x)
MLP	FC1(64) -> ReLU -> FC2(16)
Residual	x = mlp_out + x_residual
Project	lm_head: [16] -> [27] logits
Output	softmax -> next token prob

Parameter count breakdown

wte	27 x 16 = 432
wpe	16 x 16 = 256
Wq, Wk, Wv, Wo	4 x (16x16) = 1024
mlp_fc1	64 x 16 = 1024
mlp_fc2	16 x 64 = 1024
lm_head	27 x 16 = 432
TOTAL	~4,192 scalars

Key functions

<code>linear(x, w)</code>	Matrix-vector multiply (list comprehension)
<code>softmax(logits)</code>	exp(x-max)/sum with numerical stability
<code>rmsnorm(x)</code>	x / sqrt(mean(x^2) + eps)
<code>gpt(...)</code>	Full decoder pass, returns logits
<code>Value.backward()</code>	Reverse-topo chain rule

Training recipe

Dataset: 32K names from Karpathy's makemore
 Tokenizer: character-level (a-z + BOS = 27 tokens)
 Loss: cross-entropy (next character prediction)
 Optimizer: Adam (beta1=0.85, beta2=0.99)
 Learning rate: 0.01 with linear decay to 0
 Batch size: 8 names per step
 Steps: 20 (very small — for demo purposes)
 Init: Gaussian(mean=0, std=0.08) for all weights

Attention math (per head)

- Q = linear(x, Wq) -> split by head
- K = linear(x, Wk) -> append to cache
- V = linear(x, Wv) -> append to cache
- scores = Q . K^T / sqrt(head_dim)
- weights = softmax(scores) # causal: sees past only
- output = weights . V # weighted sum
- Concat all heads -> output projection Wo

Adam update (per parameter)

```
m = 0.85*m + 0.15*grad # momentum
v = 0.99*v + 0.01*grad^2 # velocity
m_hat = m / (1 - 0.85^t) # bias correct
v_hat = v / (1 - 0.99^t) # bias correct
param -= lr * m_hat / (sqrt(v_hat)+1e-8)
param.grad = 0 # zero grad
```

Why this file matters

Everything is scalar — no tensor abstraction hiding anything
 Same architecture as GPT-2 / GPT-3 (just tiny)
 Same autograd as PyTorch (topo sort + chain rule)
 Same Adam optimizer (moment estimates + bias correction)
 Same KV cache trick used in production LLMs
 ~200 lines = the whole story from scratch