# Token Generation & Loss (Forward)

$$\mathbf{A}_{\text{out}} \xrightarrow{\;[B,\,S,\,D]\;} \odot \xrightarrow[\;[B,\,S,\,V]\;]{\mathbf{z}_{\text{lin}}} + \xrightarrow[\;[B,\,S,\,V]\;]{\mathbf{z}_{\text{bias}}} \boxed{\text{S}} \xrightarrow[\;[B,\,S,\,V]\;]{\mathbf{P}}$$

$\odot$ input from below: $[D,\,V]$

$$\widetilde{\mathbf{W}}_{\text{lm}} = \mathbf{E}^T$$

$+$ input from below: $\text{BC}_{B,S}(\widetilde{\mathbf{b}}_{\text{lm}})$  $[V]$

$$\widetilde{\mathbf{b}}_{\text{lm}}$$

$\boxed{\text{TOP-}k} \xrightarrow[\;[B,\,S]\;]{\text{token ids}} \hat{\mathbf{Y}}_{\text{sample}}$

$\boxed{\text{ARG}} \xrightarrow[\;[B,\,S]\;]{\text{token ids}} \hat{\mathbf{Y}}_{\text{greedy}}$

$\boxed{\text{CE}} \xleftarrow[\;[B,\,S]\;]{\text{targets}} \mathbf{Y}_{\text{targets}}$

$$\mathcal{L}$$
(mean over $B, S$)

# Token Generation & Loss — Backward (Corrected)

$\mathrm{d}\widetilde{\mathbf{W}}_{\mathrm{lm}}$
$[D,V]$

$\mathbf{A}_{\mathrm{out}}$ —→ T $[B, D, S]$ →(•)←

$\mathrm{d}\mathcal{L}$

$\mathrm{d}\mathbf{A}_{\mathrm{out}}$ $[B, S, D]$

$\mathrm{d}\mathbf{A}_{\mathrm{out}}$ ←(•)← $\mathrm{d}\mathbf{Z}_{\mathrm{lin}} = \mathrm{d}\mathbf{Z}_{\mathrm{bias}}$ $[B, S, V]$ ← dS ← dCE

T

$\mathbf{P}$

$\mathbf{Y}_{\mathrm{targets}}$

$\widetilde{\mathbf{W}}_{\mathrm{lm}}$
$[D,V]$

$\Sigma_{B,S}$

$\mathrm{d}\widetilde{\mathbf{b}}_{\mathrm{lm}}$
$[V]$

| | Abbrev | Name | Type / Shape | Notes |
|---|---|---|---|---|
| **Operations (Ops)** | S | Softmax | op | Over vocab axis $V$; outputs probabilities $\mathbf{P}$. |
| | CE | Cross-Entropy | op | Usually *sparse* CE consuming label indices $\mathbf{Y}$. |
| | ARG | Argmax (greedy) | op | $\mathrm{argmax}_V$ to get token ids (no gradient). |
| | TOP-$k$ | Top-$k$ / sampling | op | Optional decoding path (or nucleus sampling); no gradient. |
| | T | Transpose | op | E.g., $\widetilde{\mathbf{W}}_{\mathrm{lm}}^T \in \mathbb{R}^{V \times D}$. |
| | $\mathrm{BC}_{B,S}(\cdot)$ | Broadcast | op | Expand $[V] \rightarrow [B, S, V]$ for bias add. |
| | dS | Softmax backward | op | The output is $\mathrm{d}\mathbf{Z}_{\mathrm{bias}} = \mathbf{P} - \mathrm{onehot}(\mathbf{Y})$ (with CE). |
| | dAddB | Addition (Bias) backward | op | Passes $\mathrm{d}\mathbf{Z}_{\mathrm{bias}}$ to $\mathrm{d}\mathbf{Z}_{\mathrm{lin}}$ and $\sum_{B,S}$. |
| | $\sum_{B,S}$ | Summation | op | Sums $\mathrm{d}\mathbf{Z}_{\mathrm{bias}}$ over axes $B$ and $S$ to get $\mathrm{d}\widetilde{\mathbf{b}}_{\mathrm{lm}}$. |


## Data Tensors (Values)

| Symbol | Name | Shape | Notes |
|---|---|---|---|
| $\mathbf{A}_{\mathrm{out}}$ | Transformer output (hidden) | $[B, S, D]$ | Final hidden from the Transformer block(s). |
| $\widetilde{\mathbf{W}}_{\mathrm{lm}}$ | LM head weight (tied) | $[D, V]$ | Typically tied to $\mathbf{E}^T$. |
| $\widetilde{\mathbf{b}}_{\mathrm{lm}}$ | LM head bias | $[V]$ | Broadcast-added over $[B, S, V]$. |
| $\mathbf{Z}_{\mathrm{lin}}$ | Logits (linear output) | $[B, S, V]$ | $\mathbf{Z}_{\mathrm{lin}} = \mathbf{A}_{\mathrm{out}}\widetilde{\mathbf{W}}_{\mathrm{lm}}$. |
| $\mathbf{Z}_{\mathrm{bias}}$ | Logits (final/Softmax input) | $[B, S, V]$ | $\mathbf{Z}_{\mathrm{bias}} = \mathbf{Z}_{\mathrm{lin}} + \widetilde{\mathbf{b}}_{\mathrm{lm}}$. |
| $\mathbf{P}$ | Probabilities | $[B, S, V]$ | $\mathbf{P} = \mathrm{softmax}(\mathbf{Z}_{\mathrm{bias}})$. |
| $\mathbf{Y}$ | Target token ids | $[B, S]$ | Ground-truth indices (sparse labels). |
| $\mathcal{L}$ | Loss | scalar or $[B, S]$ | Typically mean over $B, S$. |
| $\mathrm{d}\mathcal{L}$ | Loss gradient | scalar-grad | Starting signal for backward pass. |
| $\mathrm{d}\mathbf{Z}_{\mathrm{bias}}$ | Final Logits gradient | $[B, S, V]$ | From CE+Softmax: $\mathbf{P} - \mathrm{onehot}(\mathbf{Y})$. |
| $\mathrm{d}\mathbf{Z}_{\mathrm{lin}}$ | Linear output grad | $[B, S, V]$ | Same as $\mathrm{d}\mathbf{Z}_{\mathrm{bias}}$ (input to $\mathbf{Z}_{\mathrm{lin}}$ matmul). |
| $\mathrm{d}\widetilde{\mathbf{W}}_{\mathrm{lm}}$ | LM weight grad | $[D, V]$ | $= \mathbf{A}_{\mathrm{out}}^T\mathrm{d}\mathbf{Z}_{\mathrm{lin}}$. |
| $\mathrm{d}\widetilde{\mathbf{b}}_{\mathrm{lm}}$ | LM bias grad | $[V]$ | $= \sum_{B,S}(\mathrm{d}\mathbf{Z}_{\mathrm{bias}})$. |
| $\mathrm{d}\mathbf{A}_{\mathrm{out}}$ | Hidden grad | $[B, S, D]$ | $= \mathrm{d}\mathbf{Z}_{\mathrm{lin}}\widetilde{\mathbf{W}}_{\mathrm{lm}}^T$. |

**Shapes:** $B$=batch, $S$=sequence length, $D$=hidden dim, $V$=vocab size.