



This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

SPRING AI

GENERATIVE ARTIFICIAL INTELLIGENCE CON JAVA

Simone Scannapieco

Corso avanzato per Venis S.p.A, Venezia, Italia

Novembre 2025

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and extend across the width of the page, providing a template for writing or drawing. There are no margins, text, or other markings on the paper.

- ❓ ... ma come usarli per *task* specifici o con conoscenza che a loro manca?!

[illegible]



This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

- ➔ Creare una architettura neurale da zero...
- ➔ ... oppure scegliere una architettura in letteratura (per i meno sadici)
- ➔ Addestramento da zero (a partire da pesi e *bias random*)

- ➔ Sfruttare una rete neurale già addestrata su un altro insieme di dati di addestramento
- ➔ Modificare solo alcuni strati (solitamente gli ultimi) per addestrare la rete per i propri scopi

<i>Computer Vision</i>	<i>Full learning</i>	<i>Transfer learning</i>
Numero dati addestramento	10^3-10^6	10^2
Computazione	Intensiva (GPU)	Media (CPU-GPU)
Tempo di addestramento	Giorni-settimane	Ore-giorni
Accuratezza del modello	Alta	Variabile

Note

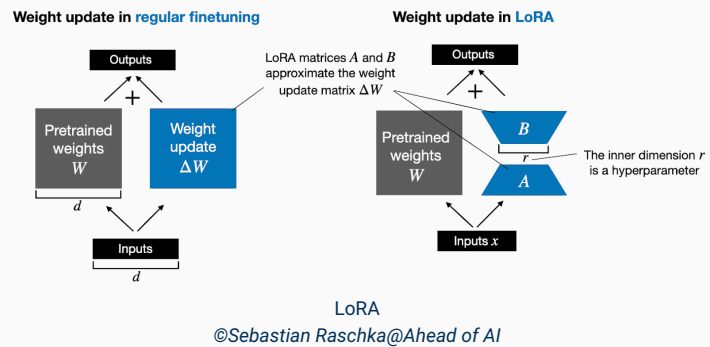
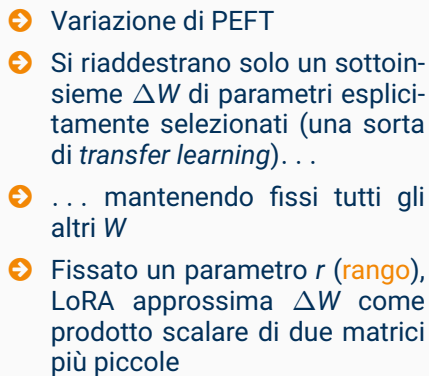
[illegible]

-  Simone Scannapieco

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 20 lines visible. The paper has a slight shadow on the right side, suggesting it's resting on a surface.

-  Simone Scannapieco

[illegible]



Note

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

Il *prodotto scalare* di matrice $(n \times r)$ $A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1r} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nr} \end{bmatrix}$ e matrice $(r \times m)$

$$B = \begin{bmatrix} b_{11} & b_{12} & \dots & b_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ b_{r1} & b_{r2} & \dots & b_{rm} \end{bmatrix} \quad \text{é la matrice } (n \times m)$$

$$A \cdot B = \begin{bmatrix} a_{11} * b_{11} + \dots + a_{1r} * b_{r1} & \dots & a_{11} * b_{1m} + \dots + a_{1r} * b_{rm} \\ \vdots & \ddots & \vdots \\ a_{n1} * b_{11} + \dots + a_{nr} * b_{r1} & \dots & a_{n1} * b_{1m} + \dots + a_{nr} * b_{rm} \end{bmatrix}$$

- ➔ In pratica, fissato r , LoRA computa A e B tale per cui $\Delta W = A \cdot B$
- ➔ Ma perché é così potente?!

Note

[illegible]

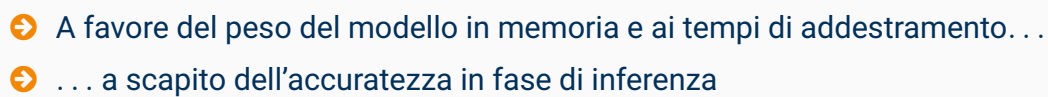
 Simone Scannapieco
 Spring AI - Corso avanzato
 Venis S.p.A, Venezia, IT
 11 / 19

[illegible]

- ➔ Quantizzazione a `float16` e `bfloat16` usati maggiormente per addestramento
- ➔ `bfloat16` generalmente preferito a `float16`
- ➔ Addestramento a `float32` riservato alle *big companies*
- ➔ Quantizzazioni inferiori disponibili (`int8`, `int4`), ma consigliate per inferenza

Note

This image shows a blank sheet of white paper with horizontal ruling lines. The lines are evenly spaced and run across the width of the page. There are no margins, text, or other markings on the paper.

[illegible]

- | Formato | Riduzione memoria | Uso principale | Accuratezza |
|---------|-------------------|-----------------|-------------|
| int8 | ~50% | Inferenza | Alta |
| int4 | ~75% | Inferenza | Media-Alta |
| GPTQ | ~75% | Inferenza (GPU) | Alta |
| GGUF | 50-80% | Inferenza (CPU) | Variabile |

[illegible]



- 15 / 19

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page. There are approximately 20 lines visible. The paper has a slight shadow on its right side, suggesting it's resting on a surface.

-  Simone Scannapieco

Note

This image shows a single sheet of white paper with horizontal blue or grey ruling lines. The lines are evenly spaced and run across the width of the page, typical of notebook paper. There are no margins, text, or other markings on the page.

- The screenshot displays the 'LLM Inference: VRAM & Performance Calculator' interface. The left sidebar contains configuration options: 'Inference' (selected) and 'Fine-tuning'. Under 'Inference', the 'Select Model' dropdown is set to 'DeepSeek-R1.5B'. The 'Inference Quantization' section is set to 'FP16'. The 'KV Cache Quantization' section is set to 'FP16/FP32 (Default)'. The 'Hardware Configuration' section shows 'RTX 3090 (24GB)' for VRAM and '8' for GPU count. The 'Batch Size' is set to '1' and 'Log Scale' is checked. The 'Segment Length' is set to '1'. The 'Concurrent Users' is set to '1'. The right panel, titled 'Performance & Memory Results', shows a green circular progress indicator at 62.7% (Good). The performance is 'MODERATE' with a score of '7.52 GB'. Below this, a list of metrics is shown: 'Generation Speed - 38 tokens/s', 'KV Cache Memory - 10.5 GB', 'Time to First Token - 0.23ms', 'Total Throughput - 38 tokens/s', and 'Profile: Optimized for Lowest Latency'. The selected model is 'DeepSeek DeepSeek-R1.5B' with links for 'Weights', 'KV Cache', 'Model Details', and 'Performance Benchmarking'. At the bottom, it says 'Model: Inference / Batch: 1'.

17 / 19

Note

[illegible]

- ➔ Adattamento stile, tono, formato *output*
- ➔ Comportamenti specifici o *task* strutturati
- ➔ *Dataset* 1K–100K esempi, risorse *hardware* limitate
- ➔ Necessità di gestire multipli adattatori per *task* diversi

[illegible]