

## Level Up Coding

★ Member-only story

# Stepping into the Future of 3D Vision with VGGT

How Transformers Are Redefining 3D Reconstruction — Faster, Smarter, and Geometry-Free



Cristian Leo · Following

Published in Level Up Coding · 26 min read · 1 hour ago



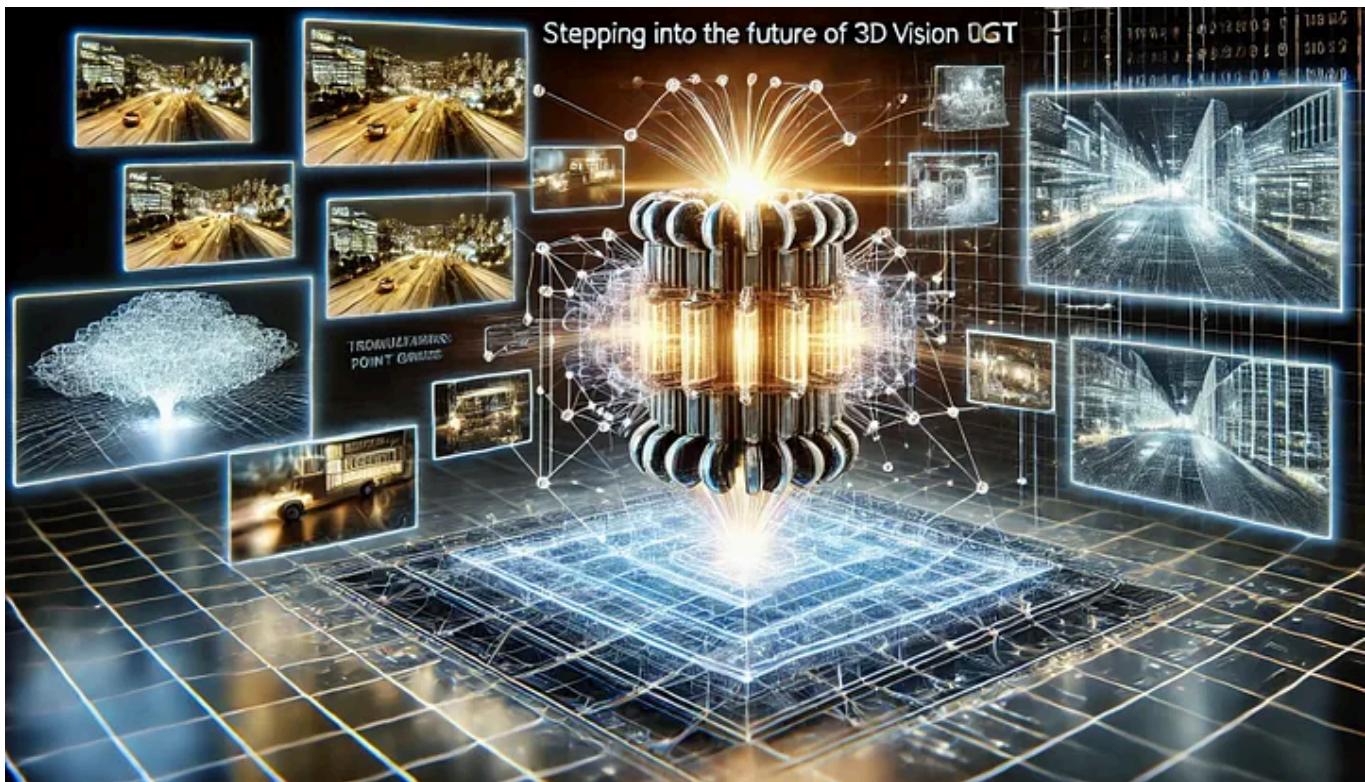


Image by DALL-E

**3D reconstruction.** It's a term that might sound like it belongs in a sci-fi movie, but it's actually becoming increasingly part of our daily lives. Think about it: from the maps guiding our self-driving cars to the immersive worlds of virtual reality, understanding and recreating the 3D structure of our environment is becoming less of a futuristic dream and more of an essential technology. For years, achieving this has felt like a complex puzzle, pieced together painstakingly using methods rooted in visual geometry. These traditional approaches, often rely on clever techniques like **Bundle Adjustment** [[Richard Hartley and Andrew Zisserman. Multiple View Geometry in Computer Vision](#)], have been the pillars, meticulously optimizing and refining 3D models from images. And while incredibly effective in many scenarios, they can also be quite demanding — computationally intensive, sometimes intricate to set up, and often requiring a good deal of post-processing to really shine.

Then, along came machine learning, initially playing a supporting role. It helped with tasks that geometry alone found tricky, like figuring out if features in two images actually matched or predicting depth from a single photo. Gradually, this integration became tighter, leading to fascinating end-to-end systems, where machine learning and visual geometry danced together, like in the impressive VGGSFm approach [[Wang et al. VGGSFm: visual geometry grounded deep structure from motion](#)]. Even with these advancements, the core of 3D reconstruction still often leaned heavily on those visual geometry principles, which, while reliable, could sometimes feel like they were adding layers of complexity and computational cost.

But what if, just what if, we could bypass much of this intricate geometric post-processing altogether? As neural networks grow in power and sophistication, a question arises: could we directly teach a network to understand and reconstruct 3D scenes, almost like magic? It's a bold idea, and recent models like DUSt3R [[Wang, et al. DUSt3R: Geometric 3D vision made easy](#)] and its successors have started to hint at this possibility, showing promising results by directly predicting 3D structures. However, these initial steps often came with their own set of constraints, perhaps limited to processing pairs of images or still needing some form of post-processing to truly scale.

This brings us to **VGGT: the Visual Geometry Grounded Transformer**. Imagine a neural network that, in a single, swift forward pass, can take in not just one or two, but potentially *hundreds* of images of a scene and instantly predict a comprehensive set of 3D attributes: camera positions, depth maps, point clouds, even track points across images. And remarkably, according to the research paper introducing VGGT, it often outperforms those optimization-based methods, and does so *without* needing further processing after its initial prediction! This is a rather substantial departure

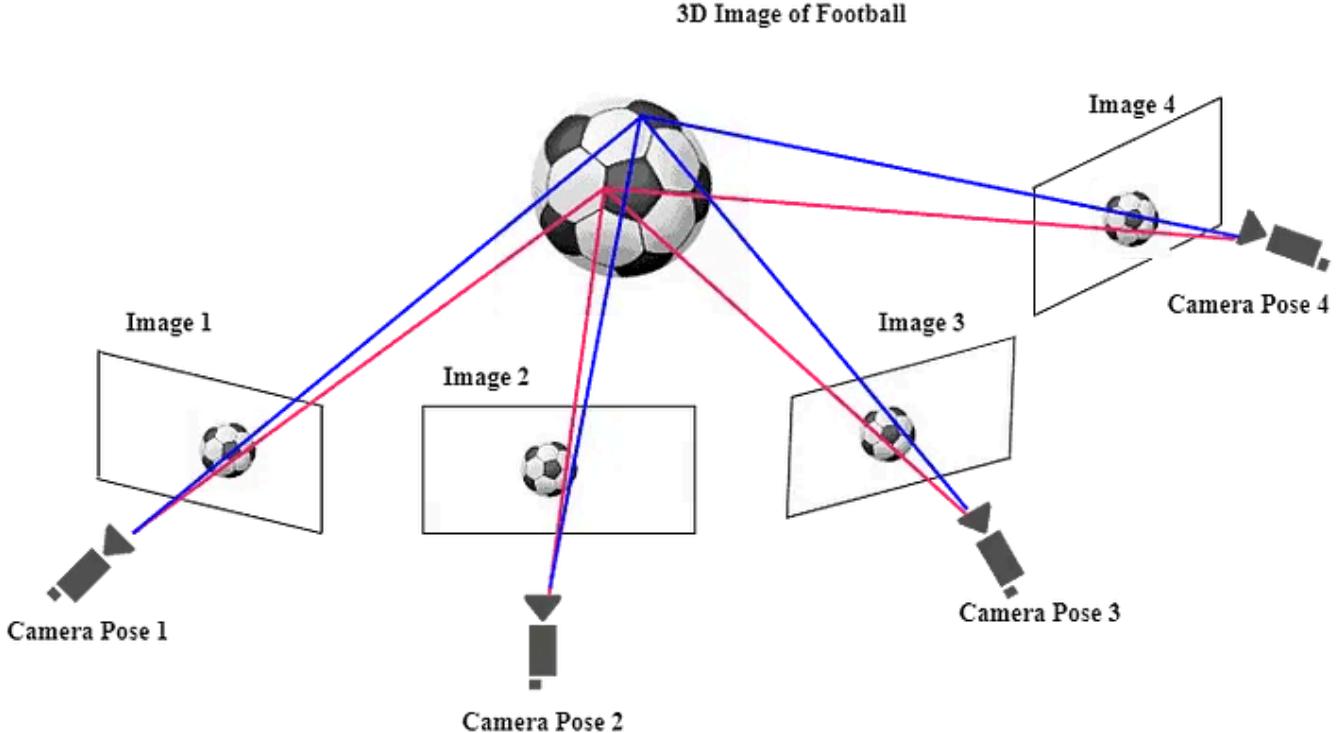
from the previous generation of models. It's like moving from carefully hand-drawing a map to instantly generating a detailed satellite view.

Is this the dawn of a new era in 3D computer vision? Is VGGT truly a step towards a future where 3D reconstruction becomes dramatically simpler, faster, and more accessible? In this article, we're going to delve deep into VGGT, unpack its architecture, and explore the math that makes it tick. We'll trace the evolution of 3D reconstruction, from its geometric roots to this exciting neural network approach. And, of course, we'll take a balanced look, considering not just the impressive potential but also the practical implications and limitations of this fascinating technology. Join me as we explore the world of VGGT and ponder what it might mean for the future of how computers see — and understand — our three-dimensional world.

## From Geometry to Neural Networks in 3D Reconstruction

### Tracing the Evolution of 3D Scene Understanding

To truly appreciate what VGGT brings to the table, it's essential to take a little trip back in time, to understand the path that 3D reconstruction has traversed. For decades, before the current wave of deep learning, the field was dominated by what we often call "classical" methods, deeply rooted in the principles of visual geometry. If you were to picture the quintessential tool in this era, it would undoubtedly be **Bundle Adjustment (BA)** [[Richard Hartley and Andrew Zisserman. Multiple View Geometry in Computer Vision](#)].



Reconstructions from multiple images and camera views — Credits by [Dhiraj Kumar](#)

Think of Bundle Adjustment as the master craftsman of 3D reconstruction. It's an iterative optimization technique, almost like a sculptor meticulously refining a clay model. Imagine trying to reconstruct a scene from a set of photographs taken from different angles. BA essentially works by simultaneously adjusting ("bundling") the camera parameters — where each camera was, its orientation — and the 3D coordinates of points in the scene. It does this by trying to minimize the discrepancies between where the 3D points *should* project in the images (based on the current estimates) and where they are *actually* observed. For a long time, BA and techniques built around it were the gold standard, powering applications from mapping entire cities to creating detailed 3D models of historical artifacts.

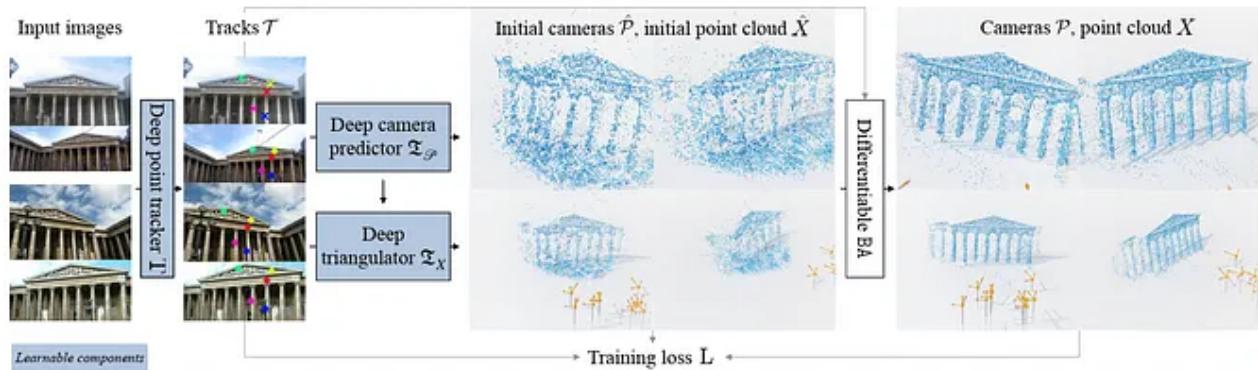
However, classical visual geometry methods, while powerful, also came with their own set of challenges. The iterative nature of Bundle Adjustment, for instance, while enabling refinement, could also be computationally quite demanding, especially for large scenes or with many images. Furthermore,

these methods often relied on robust initial estimates and could sometimes struggle in scenarios with poor texture, significant occlusions, or when dealing with vast datasets. Setting up a robust SfM (Structure from Motion) pipeline, while well-established, could also involve a degree of expertise and careful parameter tuning.

Then, the landscape began to shift. Machine learning, with its remarkable ability to learn complex patterns from data, started to step onto the scene, initially in a somewhat complementary role. It wasn't about replacing geometry entirely, at least not at first, but rather about augmenting it, tackling tasks where geometry alone faced limitations. Consider feature matching – figuring out if a point in one image corresponds to the same point in another. Classical methods existed, but machine learning, particularly with the advent of trainable feature descriptors, offered the potential for more robust and discriminative matching, especially in challenging conditions. Similarly, monocular depth prediction, estimating depth from a single image, was another area where machine learning started to shine, learning intricate cues from image appearance that pure geometry couldn't easily capture.

This integration of machine learning and visual geometry became increasingly sophisticated. Researchers started to explore “differentiable SfM” pipelines, where machine learning modules were woven into the traditional SfM framework, often with differentiable Bundle Adjustment at the core. Methods like VGGsFm [[Wang et al. VGGsFm: visual geometry grounded deep structure from motion](#)] emerged, representing a blend of both worlds. They leveraged the power of neural networks for feature extraction and matching, but still relied on the geometric principles of Bundle Adjustment for the final 3D refinement. It felt like the best of both

worlds, marrying the robustness and learning capacity of neural networks with the geometric rigor of classical techniques.



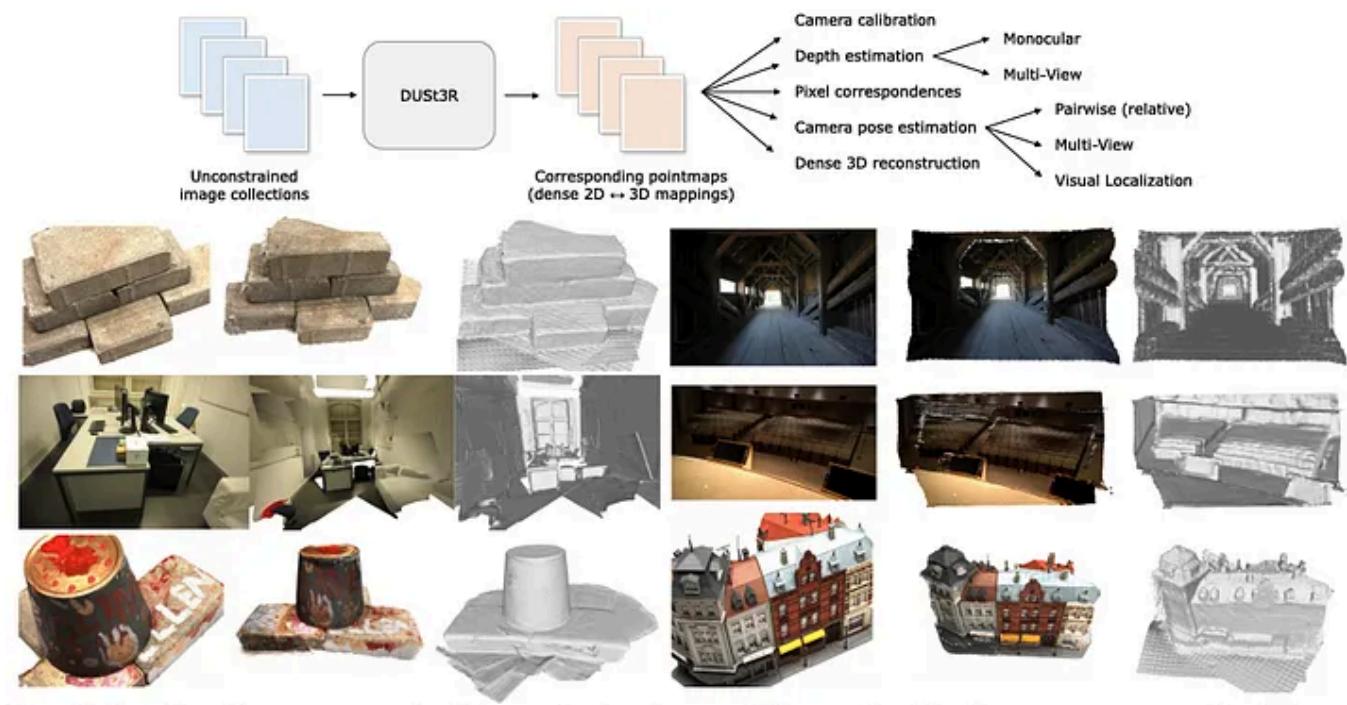
Overview of VGGSfM — Image by Wang et al. in “[VGGSf: visual geometry grounded deep structure from motion](#)”

These hybrid approaches pushed the boundaries of what was achievable, demonstrating impressive performance, particularly in challenging scenarios like large-scale phototourism reconstruction. Yet, even in these advanced hybrid systems, visual geometry, and especially Bundle Adjustment, often remained a central, if not computationally dominant, component.

But as neural networks continued their relentless march forward, growing ever larger and more capable, a truly transformative question began to be asked: could we, perhaps, finally, solve the 3D reconstruction puzzle almost *entirely* with neural networks? Could we train a network to directly infer 3D attributes from images, bypassing much of the iterative geometric post-processing?

Recent years have seen initial forays into this neural-first approach. Models like DUSt3R [[Wang, et al. DUSt3R: Geometric 3D vision made easy](#)] and MASt3R [[Leroy, et al. Grounding image matching in 3d with mast3r](#)] have

shown promising results by directly estimating aligned dense point clouds, essentially learning to perform 3D reconstruction in a more direct, data-driven manner. These were exciting steps, suggesting that a purely neural pathway to 3D vision might be within reach. However, these early networks often still had limitations. DUST3R, for example, while innovative, was primarily designed to process pairs of images and often relied on post-processing steps like global alignment to achieve high-quality results on larger scenes. MAST3R, while enhancing upon DUST3R, still operated on pairs and similarly benefited from post-optimization.



Overview of DUST3R — Image by Wang, et al. in “DUST3R: Geometric 3D vision made easy”

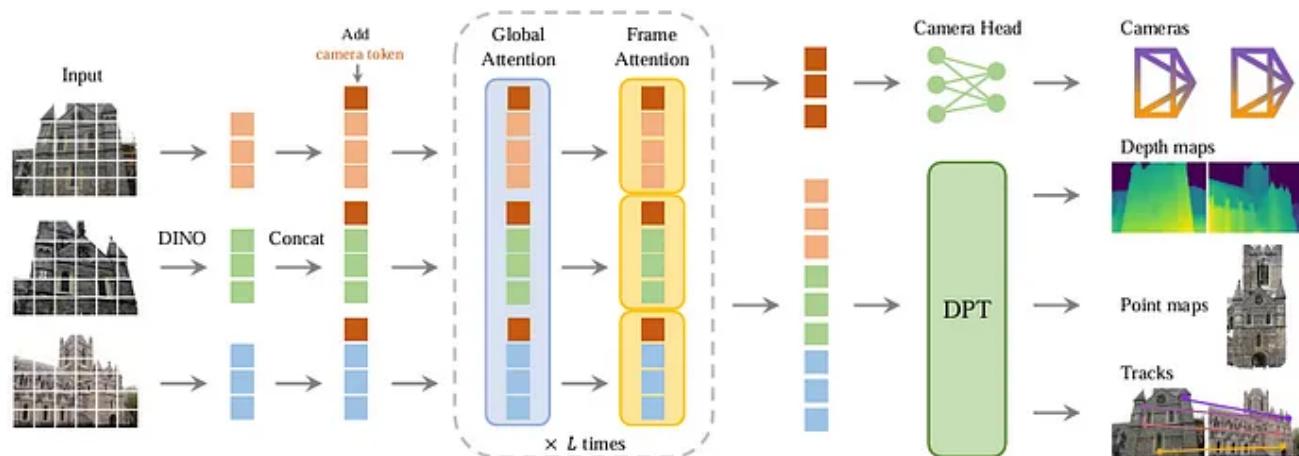
And this is where **VGGT**, the Visual Geometry Grounded Transformer, comes in. VGGT proposes a truly feed-forward neural network, based on the powerful Transformer architecture, that aims to predict a comprehensive set of 3D scene attributes — cameras, depth, point maps, tracks — in a single pass, from one, a few, or even hundreds of input images. What's particularly striking is that it reportedly achieves state-of-the-art results while often

outperforming optimization-based alternatives, even without requiring further post-processing. This is quite a claim, and it could indeed represent a substantial shift, potentially paving the way for a new generation of 3D vision systems that are not only accurate but also remarkably fast and efficient. It certainly begs the question: are we witnessing a fundamental change in how we approach 3D reconstruction, moving from a geometry-centric paradigm to a more neural, data-driven future? Let's delve deeper into the inner workings of VGGT to understand how it aims to achieve this.

## The Math Behind VGGT

### Deconstructing the 3D Inference Engine

Now, let's get to the heart of VGGT — the machinery that allows it to perform its 3D magic. To truly understand how VGGT achieves its impressive feats, we need to peek under the hood and examine the mathematical principles that underpin its architecture.



VGGT Architecture Overview — Image by Wang et al in “[VGGT:Visual Geometry Grounded Transformer](#)”

At its core, VGGT is a transformer network, a type of neural architecture that has revolutionized not only natural language processing but also,

increasingly, computer vision. If we glance at the image above, we can trace the general flow of information. The process begins with an input of a sequence of  $N$  RGB images, which we can denote as  $(I_i)_{i=1}^N$ . These images are our raw sensory data, the visual observations of a 3D scene. VGGT's task is to transform this sequence of 2D images into a set of corresponding 3D annotations for each frame. These annotations are quite comprehensive, encompassing camera parameters, depth maps, point maps, and point tracks. We can represent the output of VGGT as a set of tuples, one for each input image, like this:

$$(g_i, D_i, P_i, T_i)_{i=1}^N$$

VGGT output — Image by Author

Here,  $g_i$  represents the camera parameters,  $D_i$  is the depth map,  $P_i$  the point map, and  $T_i$  the tracking features for the  $i$ -th image.

So, how does VGGT perform this transformation? Let's start with the backbone of the network — the transformer itself and its attention mechanism.

## The Transformer Backbone: Alternating Attention (AA)

Like many modern vision models, VGGT starts by “patchifying” the input images. Think of it as dividing each image into a grid of smaller tiles, or “tokens,” each carrying local visual information. In VGGT, this patchification is achieved using DINO [[Oquab et al. DINoV2: Learning Robust Visual Features without Supervision](#)], a powerful self-supervised vision model. DINO essentially converts each input image  $I$  into a set of  $K$  tokens, which we can denote as

$$t_I \in R^{KC}$$

Image as a set of tokens K — Image by Author

Here,  $K$  represents the number of tokens, and  $C$  is the dimensionality of each token feature vector. These tokens are essentially the fundamental units of visual information that the transformer will process. The combined set of tokens from all input frames

$$t^I = \bigcup_{i=1}^N \{t_{I_i}\}$$

Combined set of tokens from all input frames — Image by Author

then becomes the input to the main network.

Now, for the heart of VGGT: the **Alternating Attention (AA)** mechanism. This is a twist on the standard transformer design, tailored for multi-image 3D reconstruction. The key idea is to make the transformer attend to the input tokens in two distinct, alternating ways: **frame-wise** and **global**.

**Frame-wise self-attention** is, in essence, about focusing on each image individually. Imagine you're looking at a single photograph. Frame-wise attention allows the network to analyze the relationships between different parts of *that same image*. It's like asking, "How do different regions within this photograph relate to each other?". Mathematically, within each frame  $i$ , the self-attention mechanism calculates attention weights based on the tokens  $\{t_{I_i}\}$  of that frame only, allowing information to flow and integrate within the spatial context of a single view.

In contrast, **global self-attention** takes a broader perspective. It attends to *all* the tokens from *all* input frames,  $t_I$ , jointly. This is where the magic of multi-view processing really starts to happen. Global attention allows the network to establish relationships and dependencies *across different images*. It's like asking, "How do features in this photograph relate to features in *all* the other photographs of the same scene?". This is crucial for tasks like multi-view stereo and point tracking, where understanding correspondences and consistencies across views is paramount.

VGGT alternates between these frame-wise and global self-attention layers, typically for  $L=24$  layers in total. This alternating strategy is quite insightful. It allows the network to strike a balance. Frame-wise attention helps in processing and normalizing information within each image independently, potentially dealing with variations in lighting or viewpoint within a single view. Global attention, on the other hand, facilitates the integration of information across multiple images, enabling the network to reason about 3D geometry from multiple perspectives and establish crucial inter-view relationships. It's a bit like having both local experts (frame-wise attention) focusing on individual images and a global strategist (global attention) orchestrating the information from all views to build a coherent 3D understanding.

## Camera Parameter Prediction

Let's move on to how VGGT predicts the camera parameters. As we mentioned earlier, VGGT uses a specific parametrization for the camera pose, denoted as  $g=[q,t,f]$ . Let's break this down:

- $q \in \mathbb{R}^4$ : This is the **rotation quaternion**, a compact and elegant way to represent 3D rotations. Quaternions, while perhaps less intuitive than Euler angles at first glance, avoid issues like gimbal lock and are often

preferred in computer graphics and robotics for representing orientations.

- $t \in R^3$ : This is the **translation vector**, representing the camera's position in 3D space.
- $f \in R^2$ : This represents the **field of view** in two dimensions, capturing the intrinsic camera parameter related to the lens's angle of view.

VGGT predicts these camera parameters using a dedicated **Camera Head**. This head takes the output camera tokens (which are special tokens appended to the image tokens) from the transformer and processes them through four additional self-attention layers. These layers further refine the camera pose representation by allowing the camera tokens to interact and refine their understanding based on the global scene context learned by the transformer. Finally, a linear layer at the end of the Camera Head maps these refined tokens to the predicted camera parameters.

## Depth Map and Point Map Prediction

VGGT also needs to predict dense outputs: depth maps and point maps. This is where the **DPT Head** comes in. The DPT Head takes the output image tokens from the transformer backbone. First, it converts these tokens into dense feature maps

$$F_i \in R^{C'' \times H \times W}$$

Dense feature maps of tokens — Image by Author

This upsampling process transforms the token-based representation back into a dense, image-like feature representation.

From these dense feature maps  $F_i$ , VGGT predicts both the **depth map**

$$D_i(y) \in R^+$$

Depth Map — Image by Author

and the **point map**

$$P_i(y) \in R^3$$

Point Map — Image by Author

for each pixel location  $y$ . This is done using a simple 3x3 convolutional layer applied to each feature map  $F_i$ . Essentially, for each pixel, the network is learning to directly regress both its depth value and its 3D coordinates in space.

A crucial point to note, as highlighted in the paper and borrowing from DUST3R [Wang, et al. DUST3R: Geometric 3D vision made easy], is that the **point maps are viewpoint invariant**. This means that the 3D points  $P_i(y)$  are defined in a fixed world coordinate system, specifically the coordinate system of the *first* camera,  $g_1$ . This choice of a common reference frame is important for multi-view consistency and for tasks like point cloud reconstruction and tracking, as it provides a unified spatial representation of the scene across different viewpoints.

## Tracking Features and Loss Functions

While VGGT itself doesn't directly output point tracks, it does predict **tracking features** through the DPT head. These dense feature maps are

designed to be fed into a separate tracking module, like CoTracker [Karaev in [CoTracker3: Simpler and Better Point Tracking by Pseudo-Labelling Real Videos](#)], to actually perform the point tracking. This modular design is quite clever. It allows VGGT to focus on learning robust and general-purpose 3D representations, while delegating the task-specific aspects of tracking to a dedicated module.

Finally, to train VGGT, the authors employ a **multi-task loss function**, a common strategy when training networks to predict multiple outputs simultaneously. The total loss  $L$  is a weighted sum of four individual loss terms:

$$L = L_{camera} + L_{depth} + L_{pmap} + \lambda L_{track}$$

Total Loss — Image by Author

Let's dissect each term:

### Camera Loss ( $L_{camera}$ )

This loss term encourages the predicted camera parameters to be close to the ground truth camera parameters. It uses the Huber loss, a robust loss function that is less sensitive to outliers than the squared error loss:

$$L_{camera} = \frac{1}{N} \sum_{i=1}^N ||\hat{g}_i - g_i||_\epsilon$$

Camera Loss — Image by Author

Here,  $||\cdot||_\epsilon$  denotes the Huber loss. This loss ensures that VGGT learns to accurately estimate the camera poses for each input image, which is

fundamental for 3D reconstruction.

### Depth Loss ( $L_{depth}$ )

This loss term guides the prediction of depth maps. Following the approach of DUStr3R [Wang, et al. DUStr3R: Geometric 3D vision made easy], it incorporates an **aleatoric uncertainty loss**. Aleatoric uncertainty, in essence, captures the inherent noise and ambiguity in the data itself. The depth loss also includes a gradient-based term, common in monocular depth estimation, to encourage sharper depth boundaries. The combined depth loss looks like this:

$$L_{depth} = \frac{1}{N} \sum_{i=1}^N \left( \left\| \frac{\hat{D}_i - D_i}{\Sigma_D^i} \right\| + \left\| \frac{\odot(\nabla \hat{D}_i - \nabla D_i)}{\Sigma_D^i} \right\| - \alpha \log \Sigma_D^i \right)$$

Depth Loss — Image by Author

Here,  $D_i$  is the ground truth depth map,  $\hat{D}_i$  is the predicted depth map,  $\Sigma_D$  is the predicted depth uncertainty map,  $\odot$  denotes element-wise multiplication, and  $\nabla$  represents the gradient operator. The first two terms penalize the discrepancy between predicted and ground truth depth and depth gradients, weighted by the uncertainty. The last term

$$-\alpha \log \Sigma_D^i$$

Aleatoric uncertainty loss — Image by Author

is the aleatoric uncertainty loss itself, encouraging the network to predict higher uncertainty in regions where depth estimation is inherently ambiguous.

## Point Map Loss ( $L_{pmap}$ )

This loss term is defined analogously to the depth loss, but for the point maps  $P_i$  and their uncertainty  $\Sigma_P$ :

$$L_{pmap} = \frac{1}{N} \sum_{i=1}^N \left( \left\| \frac{\hat{P}_i - P_i}{\Sigma_P^i} \right\| + \left\| \frac{\odot(\nabla \hat{P}_i - \nabla P_i)}{\Sigma_P^i} \right\| - \alpha \log \Sigma_P^i \right)$$

This ensures that the predicted 3D point clouds are also accurate and consistent with the ground truth, again incorporating aleatoric uncertainty.

## Tracking Loss ( $L_{track}$ )

This loss term guides the learning of tracking features. It directly supervises the predicted 2D point tracks. For each ground truth query point  $y_j$  in a query image  $I_q$ , and for each image  $I_i$ , the loss is calculated as the squared Euclidean distance between the predicted 2D correspondence and the ground truth correspondence:

$$L_{track} = \frac{1}{N_q} \sum_{j=1}^{N_q} \sum_{i=1}^N \|\hat{y}_{j,i} - y_{j,i}\|^2$$

Tracking Loss — Image by Author

Here,  $N_q$  is the number of query points. The tracking loss, weighted by a factor  $\lambda=0.05$ , encourages VGGT to learn features that are useful for tracking points consistently across multiple frames.

By jointly optimizing this multi-task loss, VGGT learns to predict a comprehensive set of 3D attributes in a coherent and consistent manner. The interplay between these different loss terms is crucial for VGGT's

performance. It's not just about predicting each 3D quantity in isolation, but about learning to predict them together, in a way that respects the inherent geometric relationships and constraints of the 3D world. This, in essence, is the mathematical foundation upon which VGGT's 3D vision capabilities are built. Now that we have a grasp of the math, let's think about how we might actually implement VGGT in code.

## Running VGGT with Python

Alright, after diving deep into the math and architecture of VGGT, you might be thinking, “Okay, this is amazing, but how do I actually *use it?*”. Well, the brilliant minds behind VGGT have made it incredibly easy for us to experiment with their model!

First things first, the VGGT team has generously shared their code and pre-trained models. This means we can skip the complex implementation and jump straight to using it.

### **GitHub - facebookresearch/vggt: [CVPR 2025] VGGT: Visual Geometry Grounded Transformer**

CVPR 2025] VGGT: Visual Geometry Grounded Transformer -  
facebookresearch/vggt

[github.com](https://github.com/facebookresearch/vggt)

Let's start by getting the code onto your computer. The easiest way is to clone their official repository from GitHub. Just open your terminal and type in:

```
git clone https://github.com/facebookresearch/vggt.git  
cd vggt
```

This will download all the necessary files to your local machine and then take you into the `vggt` directory. Next, we need to install all the required libraries that VGGT depends on. You can do this by using `pip` and the `requirements.txt` file provided in the repository:

```
python -m venv venv
source venv/bin/activate
pip install -r requirements.txt
```

This command will install libraries like `torch`, `torchvision`, `numpy`, `Pillow`, and `huggingface_hub`, which are the building blocks VGGT relies on.

Once you've got everything installed, you're ready to run VGGT! Here's a simple example of how you can load the model and run inference on some images:

```
import torch
from vggt.models.vggt import VGGT
from vggt.utils.load_fn import load_and_preprocess_images

device = "cuda" if torch.cuda.is_available() else "cpu"
# bfloat16 is supported on Ampere GPUs (Compute Capability 8.0+)
dtype = torch.bfloat16 if torch.cuda.get_device_capability()[0] >= 8 else torch.

# Initialize the model and load the pretrained weights.
# This will automatically download the model weights the first time it's run.
model = VGGT.from_pretrained("facebook/VGGT-1B").to(device)

# Load and preprocess example images (replace with your own image paths)
image_names = ["path/to/imageA.png", "path/to/imageB.png", "path/to/imageC.png"]
images = load_and_preprocess_images(image_names).to(device)

with torch.no_grad():
    with torch.cuda.amp.autocast(dtype=dtype):
```

```
# Predict attributes including cameras, depth maps, and point maps.  
predictions = model(images)
```

model = VGGT.from\_pretrained("facebook/VGGT-1B").to(device) initializes the VGGT model and loads the pre-trained weights directly from Hugging Face Hub. The first time you run this, it might take a little while as it downloads the model weights (which are around 1 billion parameters). If you ever face any issues with slow loading, you can also manually download the weights from the provided link in their documentation and load them locally.

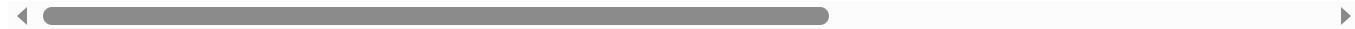
Next, you'll need to replace `["path/to/imageA.png", "path/to/imageB.png", "path/to/imageC.png"]` with the actual paths to your own images. The `load_and_preprocess_images` function takes care of loading your images and preparing them in the format VGGT expects. Finally, within the `torch.no_grad()` and `torch.cuda.amp.autocast()` context (for efficient inference), we run `predictions = model(images)`. This line performs the feed-forward pass through the VGGT network, and `predictions` will contain all the predicted 3D attributes – cameras, depth maps, point maps, and more!

For those who want a bit more control, VGGT also allows you to selectively predict specific attributes. Here's an example of how you can predict cameras, depth maps, point maps, and even get started with tracking:

```
from vggt.utils.pose_enc import pose_encoding_to_extri_intr  
from vggt.utils.geometry import unproject_depth_map_to_point_map  
  
with torch.no_grad():  
    with torch.cuda.amp.autocast(dtype=dtype):  
        images = images[None] # add batch dimension  
        aggregated_tokens_list, ps_idx = model.aggregator(images)  
        # Predict Cameras
```

```
pose_enc = model.camera_head(aggregated_tokens_list)[-1]
# Extrinsic and intrinsic matrices, following OpenCV convention (camera from
extrinsic, intrinsic = pose_encoding_to_extri_intri(pose_enc, images.shape[-1])
# Predict Depth Maps
depth_map, depth_conf = model.depth_head(aggregated_tokens_list, images, ps_
# Predict Point Maps
point_map, point_conf = model.point_head(aggregated_tokens_list, images, ps_
# Construct 3D Points from Depth Maps and Cameras
# which usually leads to more accurate 3D points than point map branch
point_map_by_unprojection = unproject_depth_map_to_point_map(depth_map.squeeze(0).float() @
extrinsic.squeeze(0).float() @
intrinsic.squeeze(0).float())
# Predict Tracks
# choose your own points to track, with shape (N, 2) for one scene
query_points = torch.FloatTensor([[100.0, 200.0],
[60.72, 259.94]]).to(device)
track_list, vis_score, conf_score = model.track_head(aggregated_tokens_list,
```

This code shows how you can access the different prediction heads of VGGT: `camera_head`, `depth_head`, `point_head`, and `track_head`. It also demonstrates how to convert the camera pose encoding into extrinsic and intrinsic matrices, and how to unproject the depth map into a point map.



The VGGT repository also provides tools for visualizing your 3D reconstructions and tracking results. They offer a Gradio web interface that lets you interactively explore the 3D scene in your browser, and a Viser 3D viewer for visualizing point clouds. For visualizing point tracks, they have a utility function to plot tracks directly on your input images. Just make sure to install the extra visualization dependencies using `pip install -r requirements_demo.txt` to use these tools. Or consider using the version hosted on HuggingFace:

**vggt - a Hugging Face Space by facebook**

VG GT (CVPR 2025)

The team has done a fantastic job of making this powerful model accessible to everyone. You can now take your own images, run them through VGGT, and start exploring the 3D world in a whole new way!

## Loss Functions and Training Loop

With the model architecture defined, the next step is to implement the loss functions we discussed in the previous section: camera loss, depth loss, point map loss, and tracking loss. You'll need to translate the LaTeX formulas into PyTorch code, using appropriate loss functions like `torch.nn.HuberLoss` and custom implementations for the aleatoric uncertainty and tracking losses.

Finally, you'll need to set up the training loop. This involves:

- **Optimizer:** AdamW is used in the paper, a common and effective optimizer for transformers.
- **Learning Rate Scheduler:** A cosine learning rate scheduler with a warmup period, as described in the paper, is typically beneficial for transformer training.
- **DataLoaders:** Use PyTorch `DataLoader` to efficiently load and batch your training data.
- **Forward Pass and Backpropagation:** In each training iteration, perform a forward pass through the VGGT model, calculate the multi-task loss, perform backpropagation, and update model weights using the optimizer.

- **Gradient Clipping:** Gradient norm clipping, with a threshold (e.g., 1.0), is often used to stabilize transformer training.
- **Mixed Precision Training:** Consider using bfloat16 mixed precision, if your hardware supports it, to speed up training and reduce memory footprint.

This, in broad strokes, outlines the key steps to implementing VGGT. It's a complex model, no doubt, but by breaking it down into modular components — feature extraction, transformer backbone, prediction heads, and loss functions — and by leveraging the power of PyTorch and transformer libraries, building a functional VGGT model, or at least experimenting with its core ideas, becomes a tangible and, I believe, a very rewarding endeavor. Of course, hyperparameter tuning, careful data preparation, and access to substantial computational resources would be essential to replicate the state-of-the-art results reported in the paper, but even a simplified implementation can provide valuable insights into the workings of this fascinating 3D vision model.

Next, let's discuss the practical implications of VGGT, its strengths, limitations, and where it might take us in the future of 3D reconstruction.

## Advantages and Considerations

Now that we've journeyed through the architecture and math of VGGT, it's time to step back and assess its practical implications. What are the real-world strengths of such a system? Where might it still fall short? And what does it all mean for the future of 3D vision? Let's weigh the advantages and considerations, trying to maintain a balanced perspective, as always.

### The Strengths of VGGT

From what we've seen in the research and the outlined architecture, VGGT appears to offer a rather compelling package of advantages:

- **Unmatched Speed and Efficiency:** Perhaps the most striking aspect of VGGT is its sheer speed. Operating in a feed-forward manner, it can process hundreds of images and predict a full suite of 3D attributes in mere seconds, as the paper claims, often under a second! This is a stark contrast to traditional optimization-based methods, which can take considerably longer, sometimes minutes or even hours, for complex scenes. In a world increasingly demanding real-time 3D perception — for applications like robotics, augmented reality, or rapid content creation — this speed advantage is, I believe, a game-changer. Imagine a self-driving car needing to reconstruct its surroundings instantaneously, or an AR application needing to understand a user's environment in the blink of an eye. VGGT's efficiency could be pivotal in enabling such real-time 3D vision.
- **State-of-the-Art Accuracy:** Speed, of course, isn't everything. But VGGT doesn't seem to compromise on accuracy either. The experimental results presented in the paper, summarized in tables 1, 2, 3, 4, and 10, suggest that VGGT achieves state-of-the-art performance across a range of 3D tasks, including camera pose estimation, depth estimation, and point cloud reconstruction. Remarkably, it often outperforms even methods that rely on computationally intensive post-optimization steps, and in some cases, even surpasses methods specialized for particular tasks. This combination of speed and accuracy is quite rare and positions VGGT as a truly competitive 3D reconstruction approach.
- **Versatile 3D Attribute Prediction:** VGGT isn't a one-trick pony. It's designed to predict a rich and diverse set of 3D scene properties simultaneously — camera parameters, depth maps, point maps, and

tracking features. This versatility, I think, is a significant strength. Instead of needing separate models for each task, VGGT provides a unified 3D representation that can be leveraged for various downstream applications. Need to estimate camera poses? VGGT's got you covered. Want a dense 3D point cloud? It's there. Interested in tracking points across frames? VGGT provides the features for that too. This multi-task capability simplifies system design and potentially allows for synergistic learning between different 3D attributes.

- **Robust Generalization:** The paper highlights VGGT's strong generalization ability, particularly on the RealEstate10K dataset, which it wasn't explicitly trained on. This suggests that VGGT, trained on a diverse and large-scale dataset, learns robust 3D representations that are not overly specialized to its training data. Generalization is a critical factor for real-world deployment, where models often encounter scenes and environments different from their training data. VGGT's apparent ability to generalize well makes it potentially more practical and widely applicable.
- **Adaptability for Downstream Tasks:** Beyond its core 3D reconstruction capabilities, VGGT's learned features seem to be highly adaptable for downstream tasks. The paper demonstrates this by fine-tuning VGGT as a feature extractor for novel view synthesis and dynamic point tracking (Section 4.6). The competitive performance achieved in these downstream tasks, even with relatively simple fine-tuning, underscores the versatility and richness of VGGT's learned 3D representations. It suggests that VGGT could serve as a powerful 3D "backbone" for a wide range of applications, much like ImageNet pre-trained networks have become foundational in 2D vision.

## Limitations and Future Directions

Despite its impressive strengths, it's important to acknowledge that VGGT, like any model, is not without limitations. The authors themselves point out a few areas for improvement:

- **Image Distortions:** VGGT, in its current form, may not gracefully handle fisheye or panoramic images, and reconstruction performance can degrade under extreme input rotations. This suggests that the model's current architecture or training regime might not fully account for very wide fields of view or significant image distortions. Addressing these limitations would likely require modifications to the input representation, attention mechanisms, or training data to explicitly handle such scenarios.
- **Non-Rigid Deformations:** While VGGT can handle scenes with minor non-rigid motions, it reportedly struggles with substantial non-rigid deformations. Think of highly dynamic scenes with significant object or scene deformation. This is a common challenge for many 3D reconstruction methods, as they often implicitly assume scene rigidity. Extending VGGT to handle more significant non-rigid scenes might require incorporating mechanisms to model and account for object or scene-level deformations, perhaps through learned deformation fields or more sophisticated temporal modeling.
- **Computational Cost:** While VGGT is fast in inference, the model itself is large (1.2 billion parameters) and training it requires significant computational resources (64 A100 GPUs for nine days, as per the paper). While inference is efficient, the training cost might be a barrier for some researchers or practitioners. Future work could explore model compression techniques, distillation, or more efficient transformer architectures to reduce the training and deployment footprint of VGGT-like models.

Looking ahead, I see several promising avenues for future research and development building upon VGGT:

- **Addressing Limitations:** Directly tackling the limitations mentioned above — handling fisheye/panoramic images, extreme rotations, and non-rigid deformations — would be valuable steps towards making VGGT more robust and broadly applicable.
- **Exploring Differentiable Bundle Adjustment:** The paper notes that combining VGGT with Bundle Adjustment post-processing further improves accuracy. Exploring tighter integration, perhaps even differentiable BA within the VGGT framework itself, could be an interesting direction to push accuracy even further, potentially bridging the gap between neural speed and geometric refinement.
- **Unsupervised and Self-Supervised Training:** VGGT is trained on a massive dataset with 3D annotations. Exploring unsupervised or self-supervised training strategies could reduce reliance on labeled 3D data, which can be expensive and time-consuming to acquire. Differentiable rendering, contrastive learning, or geometric consistency losses could be promising avenues to explore for unsupervised VGGT training.
- **Real-time and Edge Deployment:** VGGT's efficiency makes it a promising candidate for real-time and edge deployment. Further optimization, model compression, and potentially specialized hardware acceleration could pave the way for deploying VGGT on mobile devices, robots, and other resource-constrained platforms, bringing neural 3D vision to a wider range of applications.

VGGT, represents a significant stride forward in neural 3D reconstruction. It offers a compelling combination of speed, accuracy, versatility, and generalization ability. While it's not a perfect system yet, and there's

certainly room for improvement, its strengths are undeniable, and its potential to reshape the landscape of 3D computer vision seems substantial. It's a testament to the power of transformers and large-scale learning, and I believe it opens up exciting new possibilities for how machines can perceive and interact with our 3D world.

## Conclusion

VGGT, the Visual Geometry Grounded Transformer, stands out as more than just another incremental improvement in a rapidly advancing field. It feels like a meaningful shift in paradigm, a bold step towards a future where 3D vision is not only more accurate but also fundamentally faster, more efficient, and perhaps, more accessible.

Let's recap the essence of VGGT. At its heart, it's a feed-forward neural network, a type of model that, in a single sweep of computation, can ingest a set of images and output a comprehensive 3D scene understanding. It predicts not just one or two, but a whole array of crucial 3D properties — camera poses, depth maps, detailed point clouds, and even features for tracking points across views. And it does so, remarkably, in a fraction of the time taken by traditional, optimization-heavy methods, often achieving comparable or even superior accuracy in the process.

This speed advantage is perhaps VGGT's most transformative contribution. It opens up the door to real-time 3D vision in scenarios where latency is paramount — think of autonomous robots navigating complex environments, augmented reality applications seamlessly blending virtual and real worlds, or interactive 3D content creation tools that respond instantly to user input. In these domains, the ability to reconstruct a 3D scene in milliseconds, rather than seconds or minutes, can be the difference

between a system that's merely functional and one that's truly transformative.

Of course, like any technology, VGGT has its limitations. It's not yet a perfect solution for every 3D vision challenge. Its current architecture may struggle with highly distorted images or scenes with substantial non-rigid

Open in app ↗



Search



Write



Yet, looking at the bigger picture, VGGT feels like a significant step in a larger trend. It's part of a growing wave of neural-first approaches that are challenging the traditional geometry-centric paradigm in 3D computer vision. It suggests that, with the right architectures and training data, neural networks can learn to directly reason about 3D geometry in a way that's not only accurate but also remarkably efficient, potentially bypassing the need for complex and time-consuming post-processing.

Is VGGT the ultimate answer to all 3D reconstruction problems? Perhaps not yet. But it's certainly a compelling glimpse into a future where 3D vision is democratized, becoming faster, more accessible, and more deeply integrated into our everyday technologies. And that future, I believe, is one worth striving for.

## References

- Richard Hartley and Andrew Zisserman.  
*Multiple View Geometry in Computer Vision*, 2nd Edition. Cambridge University Press, 2003.  
ISBN: 9780521540513.

- Wang, Y., Zhuang, Y., Yu, C., Liu, S., Tang, M., Dai, Y., & Shen, X. (2023).  
**VGGSFm: Visual Geometry Grounded Structure-from-Motion.**  
*arXiv preprint arXiv:2312.02935.*  
<https://arxiv.org/abs/2312.02935>
- Wang, Y., Tang, M., Yu, C., Liu, S., & Shen, X. (2023).  
**DUSt3R: Geometric 3D Vision Made Easy.**  
*arXiv preprint arXiv:2305.06462.*  
<https://arxiv.org/abs/2305.06462>
- Leroy, V., Zaffar, M., Yi, K. M., & Fua, P. (2023).  
**MASt3R: Grounding Image Matching in 3D with MASt3R.**  
*arXiv preprint arXiv:2306.07379.*  
<https://arxiv.org/abs/2306.07379>
- Karaev, D., Carion, N., & Dehghani, M. (2023).  
**CoTracker: Simpler and Better Point Tracking by Pseudo-Labelling Real Videos.**  
*arXiv preprint arXiv:2307.08263.*  
<https://arxiv.org/abs/2307.08263>
- Oquab, M., Darzet, T., Moutakanni, T., Fernandez, P., Haziza, D., Sun, C., ... & Caron, M. (2023).  
**DINOv2: Learning Robust Visual Features without Supervision.**  
*arXiv preprint arXiv:2304.07193.*  
<https://arxiv.org/abs/2304.07193>
- Wang, Y., Yu, C., Liu, S., Tang, M., Dai, Y., & Shen, X. (2025).  
**VGGT: Visual Geometry Grounded Transformer.**  
*Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2025).*  
GitHub: <https://github.com/facebookresearch/vggt>  
Hugging Face: <https://huggingface.co/facebook/VGGT-1B>

What are your thoughts on VGGT and the future of neural 3D reconstruction? Share your comments and questions below!

If you found this article insightful, make sure to follow me, Cristian Leo, on Medium for more in-depth explorations of Data Science and Computer Vision breakthroughs. Let's learn and grow together in this exciting field!

Computer Vision

Artificial Intelligence

Mathematics

Deep Learning

Python



Published in **Level Up Coding**

Follow

221K Followers · Last published 1 hour ago

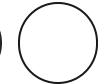
Coding tutorials and news. The developer homepage [gitconnected.com](https://gitconnected.com) && [skilled.dev](https://skilled.dev) && [levelup.dev](https://levelup.dev)



**Written by Cristian Leo**

Following

34K Followers · 13 Following



Data Scientist @ Amazon with a passion about recreating all the popular machine learning algorithm from scratch.

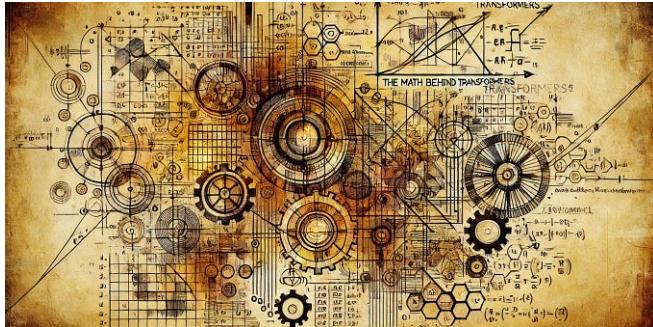
No responses yet



Alex Mylnikov

What are your thoughts?

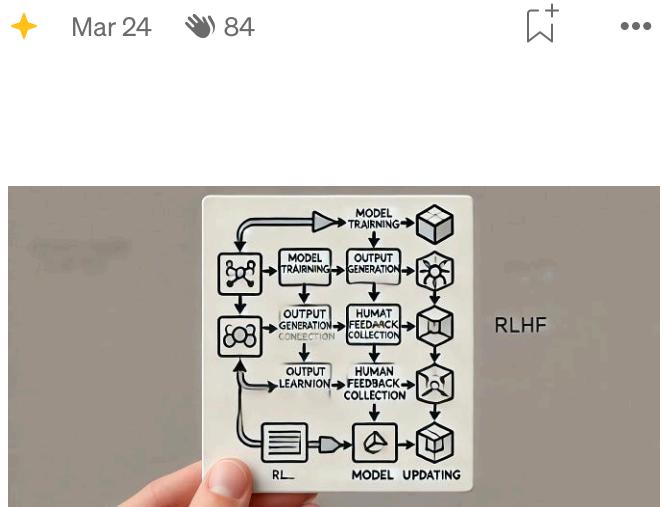
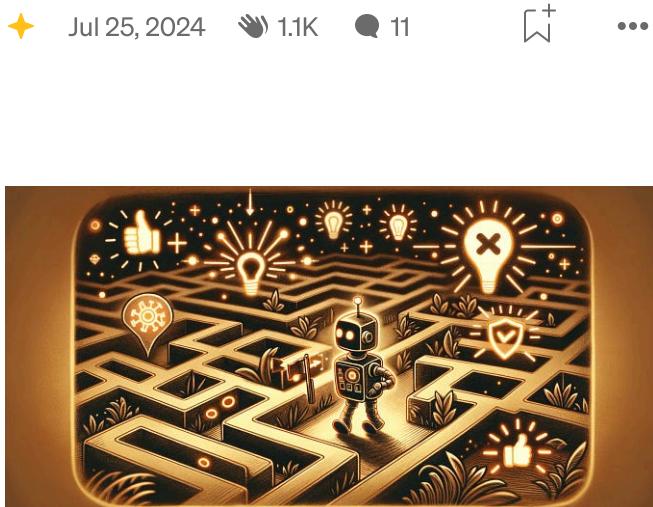
[More from Cristian Leo and Level Up Coding](#)



Cristian Leo

# The Math Behind Transformers

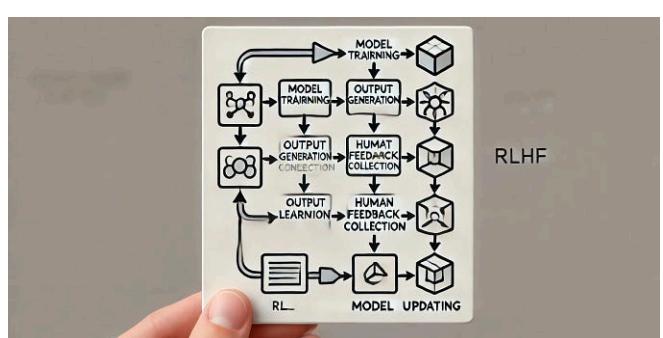
# Deep Dive into the Transformer Architecture, the key element of LLMs. Let's explore its...



In Data Science Collective by Cristian Leo

# How Qubits Are Rewriting the Rules of Computation

# From Classical Certainty to Quantum Possibility: Exploring the Science, Math, and...





In TDS Archive by Cristian Leo



In Level Up Coding by Cristian Leo

## Reinforcement Learning 101: Building a RL Agent

Decoding the Math behind Reinforcement Learning, introducing the RL Framework, an...

★ Feb 19, 2024

905

10



...

## Can We Really Trust Human Feedback?

Exploring the Limits and Potential of Reinforcement Learning from Human...

★ Mar 17

156

1

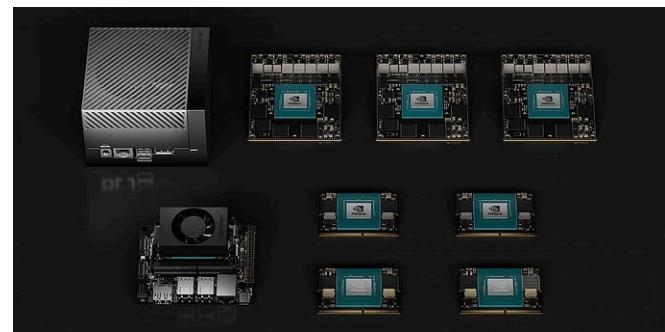
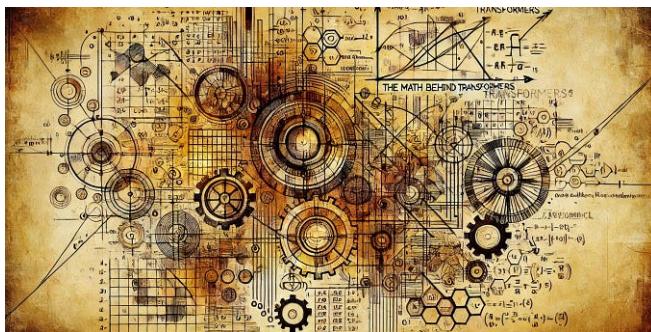


...

See all from Cristian Leo

See all from Level Up Coding

## Recommended from Medium



 Cristian Leo

## The Math Behind Transformers

Deep Dive into the Transformer Architecture, the key element of LLMs. Let's explore its...

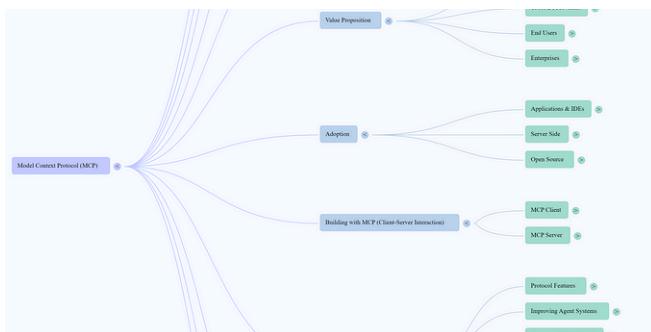
★ Jul 25, 2024 ⚡ 1.1K 🗣 11  ...

 Henry Navarro

## How to deploy YOLOv12 on NVIDIA Jetson devices? 🖨️🧠

Deploy any YOLO model on NVIDIA Jetson devices using Ultralytics and Flask.

★ Mar 17 ⚡ 11 ⚡ 11  ...



 Dr. Nimrita Koul ✨

## The Model Context Protocol (MCP) —A Complete Tutorial

Anthropic released the Model Context Protocol(MCP) in Nov. 2024.

Mar 27 ⚡ 256 🗣 1  ...



 faruk.ozelli

## YOLOv12: A New Era in Attention-Centric Real-Time Object Detecti...

Where Speed, Accuracy, and Efficiency Converge—An In-Depth Analysis of the...

★ Feb 20 ⚡ 19 ⚡ 11  ...



 Wei Lu

## Local DeepSeek-R1 671B on \$800 configurations

No matter how competitors attack DeepSeek, the V3 and R1 models are fully open-source...

 Mar 20  642  13

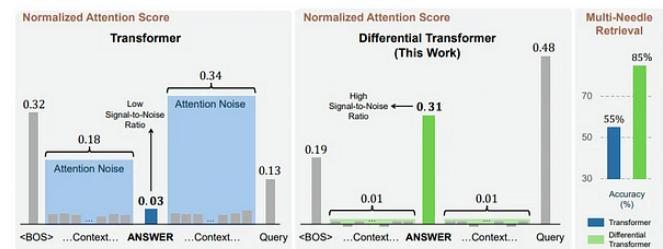


Figure 1: Transformer often over-attends to irrelevant context (i.e., attention noise). DIFF Transf. classifies attention to answer spans and cancels noise, enhancing the capability of context modeling.

 In Towards AI by Shubh Mishra

## A New Approach to Attention—Differential Transformers | Paper...

Today, we are looking into a very different approach to the transformer architecture.

Jan 28  78

[See more recommendations](#)