

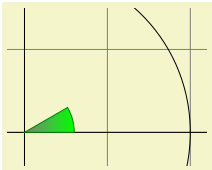
## Zero-Shot Text-Guided Graphics Program Synthesis

Jonas Belouadi Eddy Ilg Margret Keuper Hideki Tanaka Masao Utiyama  
Raj Dabre Steffen Eger Simone Ponzetto  
University of Mannheim University of Technology Nuremberg NICT



# Background

- **Graphics programming languages** offer advantages over low-level vector and raster image formats by representing visuals as **high-level** programs.

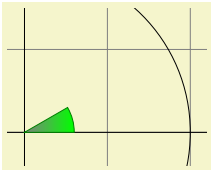


```
\begin{tikzpicture}[scale=3]
  \clip (-0.1,-0.2) rectangle (1.1,0.75);
  \draw[step=.5cm,gray,very thin] (-1.4,-1.4) grid (1.4,1.4);
  \draw (-1.5,0) -- (1.5,0);
  \draw (0,-1.5) -- (0,1.5);
  \draw (0,0) circle (1cm);
  \shadedraw[left color=gray,right color=green, draw=green!50!black]
    (0,0) -- (3mm,0mm) arc (0:30:3mm) -- cycle;
\end{tikzpicture}
```

Source: <https://tikz.dev>

# Background

- **Graphics programming languages** offer advantages over low-level vector and raster image formats by representing visuals as **high-level** programs.
  - Preserve semantics, remain human-interpretable, and allow manual editing.

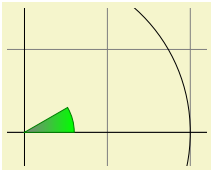


```
\begin{tikzpicture}[scale=3]
  \clip (-0.1,-0.2) rectangle (1.1,0.75);
  \draw[step=.5cm,gray,very thin] (-1.4,-1.4) grid (1.4,1.4);
  \draw (-1.5,0) -- (1.5,0);
  \draw (0,-1.5) -- (0,1.5);
  \draw (0,0) circle (1cm);
  \shadedraw[left color=gray,right color=green, draw=green!50!black]
    (0,0) -- (3mm,0mm) arc (0:30:3mm) -- cycle;
\end{tikzpicture}
```

Source: <https://tikz.dev>

# Background

- **Graphics programming languages** offer advantages over low-level vector and raster image formats by representing visuals as **high-level** programs.
  - Preserve semantics, remain human-interpretable, and allow manual editing.
- These properties are valuable to the **scientific research community**, where specialized languages like **TikZ** (and many others) are popular.

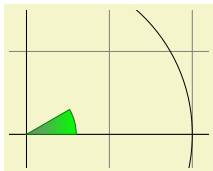


```
\begin{tikzpicture}[scale=3]
  \clip (-0.1,-0.2) rectangle (1.1,0.75);
  \draw[step=.5cm,gray,very thin] (-1.4,-1.4) grid (1.4,1.4);
  \draw (-1.5,0) -- (1.5,0);
  \draw (0,-1.5) -- (0,1.5);
  \draw (0,0) circle (1cm);
  \shadedraw[left color=gray,right color=green, draw=green!50!black]
    (0,0) -- (3mm,0mm) arc (0:30:3mm) -- cycle;
\end{tikzpicture}
```

Source: <https://tikz.dev>

# Background

- **Graphics programming languages** offer advantages over low-level vector and raster image formats by representing visuals as **high-level** programs.
  - Preserve semantics, remain human-interpretable, and allow manual editing.
- These properties are valuable to the **scientific research community**, where specialized languages like **TikZ** (and many others) are popular.
- Graphics programming languages come with a **steep learning curve** motivating **automated synthesis** approaches (e.g., with **text-guidance**).

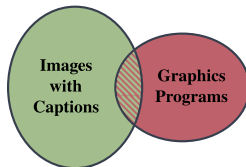


```
\begin{tikzpicture}[scale=3]
  \clip (-0.1,-0.2) rectangle (1.1,0.75);
  \draw[step=.5cm,gray,very thin] (-1.4,-1.4) grid (1.4,1.4);
  \draw (-1.5,0) -- (1.5,0);
  \draw (0,-1.5) -- (0,1.5);
  \draw (0,0) circle (1cm);
  \shadedraw[left color=gray,right color=green, draw=green!50!black]
    (0,0) -- (3mm,0mm) arc (0:30:3mm) -- cycle;
\end{tikzpicture}
```

Source: <https://tikz.dev>

# The Problem

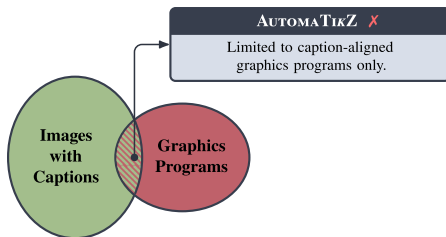
Although a **large pool** of potential training data exists for, current approaches are **restricted to subsets** based on **input modality** and **training setup**.



# The Problem

Although a **large pool** of potential training data exists for, current approaches are **restricted to subsets** based on **input modality** and **training setup**.

**AUTOMATikZ** supports text-guidance but requires graphics programs paired with captions for training (scarce) resulting in limited performance.

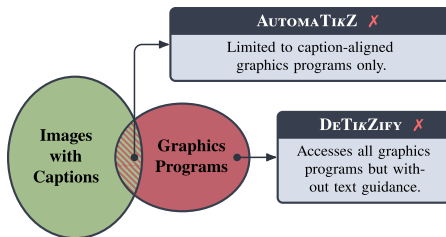


# The Problem

Although a **large pool** of potential training data exists for, current approaches are **restricted to subsets** based on **input modality** and **training setup**.

**AUTOMATikZ** supports text-guidance but requires graphics programs paired with captions for training (scarce) resulting in limited performance.

**DETIkZIFY** generates programs from images (**inverse graphics**) with self-supervised training (can use lots of training data) but creating these visual inputs remains cumbersome

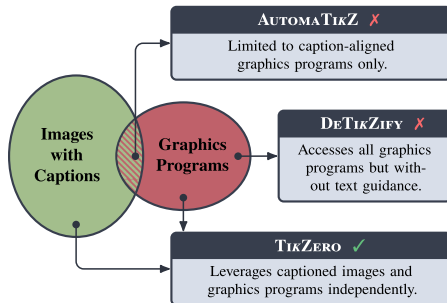




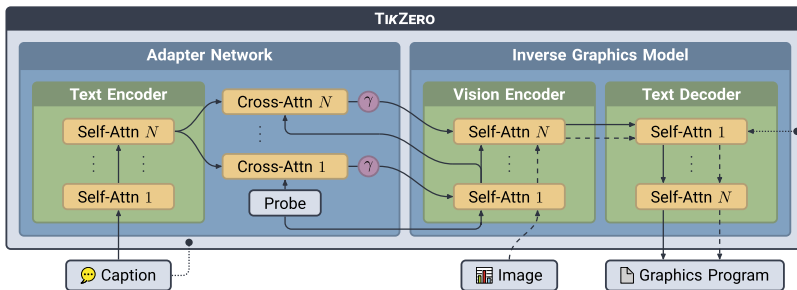
# Our TIKZERO approach

## Idea

We **decouple** graphics program generation from text understanding, enabling **independent training** on graphics programs and captioned images without requiring paired data. We call our approach **TIKZERO**.

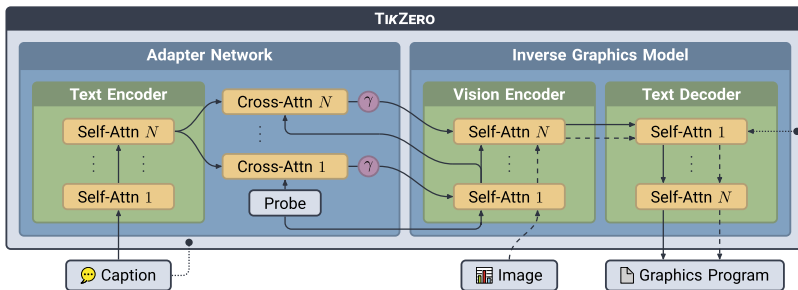


# Architecture



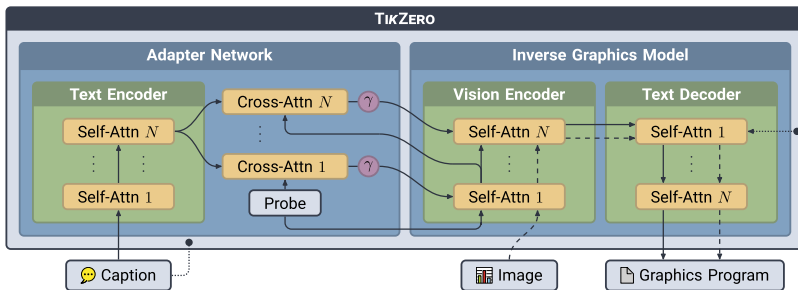
1. We train an **inverse graphics model** conditioned on **image patch embeddings** from a vision encoder (akin to DETIKZIFY).

# Architecture



1. We train an **inverse graphics model** conditioned on **image patch embeddings** from a vision encoder (akin to DETIKZIFY).
2. We then train an **adapter network** that generates **synthetic image patch embeddings** from **captions**. This adapter training relies solely on captioned images, effectively circumventing resource limitations.

# Architecture



1. We train an **inverse graphics model** conditioned on **image patch embeddings** from a vision encoder (akin to DETIKZIFY).
2. We then train an **adapter network** that generates **synthetic image patch embeddings** from **captions**. This adapter training relies solely on captioned images, effectively circumventing resource limitations.
3. Later, we combine this approach with **end-to-end fine-tuning** (TIKZERO+).

# Training Data

## Inverse Graphics Model Training Data

We systematically extract **TikZ** graphics programs from online sources. Whenever possible we also extract **captions** to support our claims. From over **450k** instances, fewer than **170k** include captions, underscoring the challenges discussed.

# Training Data

## Inverse Graphics Model Training Data

We systematically extract **TikZ** graphics programs from online sources. Whenever possible we also extract **captions** to support our claims. From over **450k** instances, fewer than **170k** include captions, underscoring the challenges discussed.

## Adapter Training Data

We leverage the **6.4 million** scientific caption-image pairs from ARXIVCAP for adapter training.

# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

**Image Similarity** DREAMSIM (DSIM); Kernel Inception Distance (KID)



# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

**Image Similarity** DREAMSIM (DSIM); Kernel Inception Distance (KID)

**Caption Similarity** CLIPSCORE (CLIP)

# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

**Image Similarity** DREAMSIM (DSIM); Kernel Inception Distance (KID)

**Caption Similarity** CLIPSCORE (CLIP)

**Code Similarity** CRYSTALBLEU (cBLEU); T<sub>E</sub>X Edit Distance (TED)

# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

**Image Similarity** DREAMSIM (DSIM); Kernel Inception Distance (KID)

**Caption Similarity** CLIPSCORE (CLIP)

**Code Similarity** CRYSTALBLEU (cBLEU); T<sub>E</sub>X Edit Distance (TED)

**Efficiency** Mean Token Efficiency (MTE)

# Automatic Evaluation

We employ the following automatic **evaluation metrics** for quantitative evaluation:

**Image Similarity** DREAMSIM (DSIM); Kernel Inception Distance (KID)

**Caption Similarity** CLIPSCORE (CLIP)

**Code Similarity** CRYSTALBLEU (CBLEU); T<sub>E</sub>X Edit Distance (TED)

**Efficiency** Mean Token Efficiency (MTE)

**Mean Similarity** average of the above (AVG)

# Results

Models	DSim <sub>↑</sub>	KID <sub>↓</sub>	CLIP <sub>↑</sub>	cBLEU <sub>↑</sub>	TED <sub>↓</sub>	MTE <sub>↑</sub>	AVG <sub>↑</sub>
IDEFICS 3 (8B)	45.475	11.426	<u>14.327</u>	0.656	63.175	69.558	66.628
AUTOMATikZ (13B)	46.033	<b>1.294</b>	3.955	0.386	<b>62.24</b>	<b>85.866</b>	63.093
AUTOMATikZ <sub>v2</sub> (VLM)	38.313	33.203	0.775	0.328	76.985	21.595	0.0
AUTOMATikZ <sub>v2</sub> (LLM)	<u>50.548</u>	<u>3.491</u>	<b>15.766</b>	<u>0.658</u>	<u>62.307</u>	81.775	<u>82.375</u>
TIKZERO	<b>52.829</b>	5.103	10.051	<b>1.603</b>	65.51	<u>82.291</u>	<b>85.599</b>

- On average, TIKZERO **outperforms** AUTOMATikZ<sub>v2</sub> (baseline trained on supervised subset of our data) and additional baselines.

## Additional Results (TIKZERO+)

Models	DSim $\uparrow$	KID $\downarrow$	CLIP $\uparrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$	AVG $\uparrow$
QWEN <sub>2.5</sub> CODER (32B)	54.473	5.493	<u>24.87</u>	0.285	59.856	<u>97.269</u>	48.593
GPT-4o	<b>56.464</b>	<u>2.844</u>	<b>31.787</b>	0.327	<b>58.511</b>	<b>97.675</b>	<u>79.019</u>
TIKZERO	52.829	5.103	10.051	<u>1.603</u>	65.51	82.291	14.658
TIKZERO+	<u>56.295</u>	<b>1.831</b>	24.177	<b>1.988</b>	<u>59.008</u>	93.058	<b>87.043</b>

- TIKZERO+ **outperforms** TIKZERO by a huge margin showing that subsequent supervised fine-tuning is beneficial.

## Additional Results (TIKZERO+)

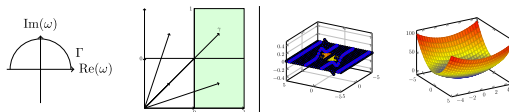
Models	DSim $\uparrow$	KID $\downarrow$	CLIP $\uparrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$	AVG $\uparrow$
QWEN <sub>2.5</sub> CODER (32B)	54.473	5.493	<u>24.87</u>	0.285	59.856	<u>97.269</u>	48.593
GPT-4o	<b>56.464</b>	<u>2.844</u>	<b>31.787</b>	0.327	<b>58.511</b>	<b>97.675</b>	<u>79.019</u>
TIKZERO	52.829	5.103	10.051	<u>1.603</u>	65.51	82.291	14.658
TIKZERO+	<u>56.295</u>	<b>1.831</b>	24.177	<b>1.988</b>	<u>59.008</u>	93.058	<b>87.043</b>

- TIKZERO+ **outperforms** TIKZERO by a huge margin showing that subsequent supervised fine-tuning is beneficial.
- It now also **outperforms** GPT-4o and other **much larger** and **commercial** models on **average** and comes close on key metrics DREAMSIM and CLIPSCORE.

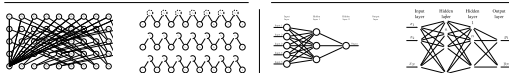
# Examples

**AUTOMATIKZ<sub>v2</sub>** ❌

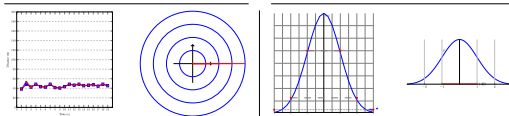
**TIKZERO+** ✅



3D contour plot of a loss function.



A multi-layer perceptron with two hidden layers.



Gaussian probability density function (blue) with markers showing one standard deviation (red).

Qualitative comparison of our **TIKZERO** approach (last two columns) and the end-to-end trained baseline **AUTOMATIKZ<sub>v2</sub>**. Our method generates outputs that **more closely** follow the given captions.



## Interested? There's more!



In our paper, we **evaluate additional baselines**, including reasoning models; uncover that end-to-end trained baselines **prioritize copying strings** over better visuals; and **evaluate the inverse graphics performance** of our model. Our code, datasets, and select models are **publicly available**.