

## A. Additional Baselines and Ablation Studies

Beyond the evaluation in Sec. 5, we test additional baselines, including reasoning models that have proven successful in program synthesis tasks [89]. We evaluate QWEN<sub>2.5</sub> CODER (14B) [75], which complements QWEN<sub>2.5</sub> CODER (32B) from Sec. 5.2, and reasoning models from the DEEPSEEK-R1 QWEN family (14B and 32B) [90]. We also evaluate TIKZERO (Base), a variant of TIKZERO (Cos) without the trainable probe and gating mechanism, to assess their contributions.

As shown in Tab. 5, TIKZERO (Cos) achieves the highest performance, surpassing TIKZERO (Base) on both DREAMSIM and CLIPSCORE metrics and in average performance. These results validate the probe and gate design. Additionally, QWEN<sub>2.5</sub> CODER (14B) performs worse than both TIKZERO (Cos) and, as expected, its 32B variant in Tab. 3. The results are consistent with our findings in Sec. 5.1 that TIKZERO (Cos) outperforms end-to-end trained baselines of comparable size. Notably, the reasoning models show the lowest overall performance, even compared to QWEN<sub>2.5</sub> CODER (14B), indicating that reasoning capabilities alone are insufficient for graphics program synthesis and more domain-specific post-training may be needed.

## B. Supplementary Comparison with DETIKZIFY

Tab. 6 shows in detail how TIKZERO’s inverse graphics model (hereafter referred to as DETIKZIFY<sub>v2</sub>) compares against DETIKZIFY<sub>DS</sub> (7b), previously the best performing DETIKZIFY model, as evaluated on the test split of DATIKZ<sub>v3</sub>. DETIKZIFY<sub>v2</sub> clearly outperforms its predecessor across all evaluated metrics. Below, we briefly outline key differences in training and inference beyond what we described in Sec. 4. For a comprehensive description of the foundation on which DETIKZIFY<sub>v2</sub> builds, we refer to Belouadi et al. [2].

**Training** Similar to DETIKZIFY, DETIKZIFY<sub>v2</sub> employs a dense layer as the modality connector between the vision encoder and text decoder. However, for pretraining this layer, we replace the METAFIG dataset [2] with the substantially larger ARXIVCAP dataset, extracting 1 million (figure, caption, OCR) triplets. During fine-tuning, we randomly substitute inputs with synthetically generated sketches to support hand-drawn inputs. To generate these sketches, we fine-tune the image-editing model ULTRAEDIT [91] on a dataset of real, human-created scientific sketches [2]. The resulting model, ULTRASKECH, achieves a congruence coefficient (CC) [92] of 0.74 with said sketches, compared to 0.72 for the previous model used with DETIKZIFY. Additionally, we generate synthetic sketches using traditional image transformations such as random displacement fields. While these sketches exhibit less diversity, they better preserve text rendering and achieve a comparable CC of 0.75. Averaging the sketch representations from both methods increases the CC to 0.82, demonstrating their complementary nature.

**Inference** DETIKZIFY implements a Monte Carlo Tree Search-based inference algorithm to iteratively refine outputs. As a reward signal  $r$ , it computes the cosine similarity  $r_{\cos} = \cos(\text{pool}(\mathbf{x}), \text{pool}(\mathbf{y}))$  between image patch embeddings  $\mathbf{x}, \mathbf{y}$  of input images and compiled outputs via a learned pooling function. Since DETIKZIFY<sub>v2</sub> fully fine-tunes the vision encoder and uses its patch embeddings directly, it cannot compute pooled embeddings in the same way. As an alternative, inspired by popular machine translation metrics [93–96], we experiment with computing the Earth Mover’s Distance (EMD) [97, 98] with image patch embeddings. Given the distance matrix  $\mathbf{D}$ , where  $D_{i,j} = \cos(x_i, y_j)$ , EMD is defined as follows:

$$\text{EMD}(\mathbf{x}, \mathbf{y}) = \frac{\sum_{i=1}^{|\mathbf{x}|} \sum_{j=1}^{|\mathbf{y}|} F_{i,j} D_{i,j}}{\sum_{i=1}^{|\mathbf{x}|} \sum_{j=1}^{|\mathbf{y}|} F_{i,j}},$$

$$\text{with } \min_{\mathbf{F} \geq 0} \sum_{i=1}^{|\mathbf{x}|} \sum_{j=1}^{|\mathbf{y}|} F_{i,j} D_{i,j} \quad (2)$$

$$\text{s.t. } \forall_{i,j} \begin{cases} \sum_{i=1}^{|\mathbf{x}|} F_{i,j} = \frac{1}{|\mathbf{y}|}, \\ \sum_{j=1}^{|\mathbf{y}|} F_{i,j} = \frac{1}{|\mathbf{x}|}. \end{cases}$$

When correlating reward scores computed as  $r_{\cos}$  from DETIKZIFY and  $r_{\text{EMD}} = \text{EMD}(x_i, y_j)$  from DETIKZIFY<sub>v2</sub> with human judgments from Belouadi et al. [2], we find that  $r_{\text{EMD}}$  enhances correlation with humans (0.456 segment-level and 0.911 system-level Spearman’s  $\rho$ ), compared to  $r_{\cos}$  (0.436 and 0.642, respectively). This demonstrates that DETIKZIFY<sub>v2</sub> not only supports the inference algorithm but improves upon DETIKZIFY’s capabilities.

## C. Supplementary Inference Details

To instruct general-purpose models to generate TikZ code, we employ a consistent prompt across all models (GPT-4o, QWEN<sub>2.5</sub> CODER (32B), and IDEFICS 3 (8B)) originally engineered by Zhang et al. [4]. For each figure, we replace the <caption> placeholder with the specific caption:

```

1 Please generate a scientific
2 figure according to the following
3 requirements: <caption>. Your output
4 should be in TikZ code. Do not include
5 any text other than the TikZ code.
```

## D. Supplementary Experimental Results

Tab. 7 presents detailed evaluation metrics scores for the low-resource training experiments discussed in Sec. 6.2. The results show a consistent degradation in performance across all metrics as both the amount of training data and the number of layers decrease, a trend effectively captured by the AVG scores also shown in Tab. 4.

## E. Annotator Demographics

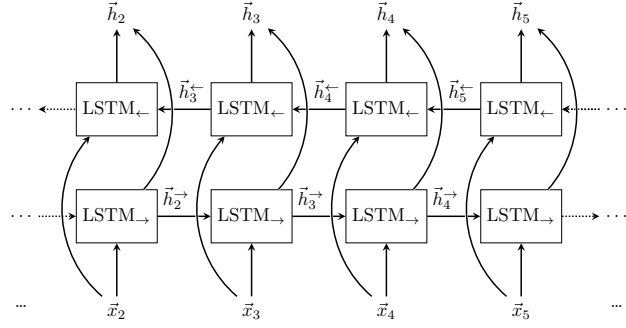
Our annotation team consists of thirteen experts with extensive research experience in Machine Learning, Natural Language Processing, or Computer Vision. The team includes one male faculty member, four female PhD students, four male PhD students, and four male researcher scientists from a research institute. We deliberately selected expert annotators based on findings by Belouadi et al. [3], which demonstrated that crowd workers often lack the necessary research background to provide reliable annotations for scientific figures. To mitigate potential biases, each annotator received the tuples and items within the tuples in randomized order.

## F. Additional Examples

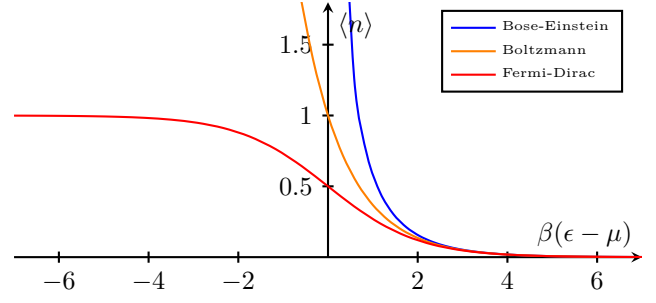
Figure 5 showcases examples<sup>5</sup> from DATIKZ<sub>v3</sub> with permissive licenses. Additionally, Tab. 8 presents randomly sampled tuples from our human evaluation with the highest and lowest rated instances highlighted. The results show that AUTOMATIKZ<sub>v2</sub> (LLM) and TIKZERO (Cos) are more frequently selected as the worst models (four and three times, respectively), while TIKZERO+ and GPT-4o are more often chosen as the best models (both three times), which aligns with our findings in Sec. 5.3. Finally, Fig. 6 illustrates example programs generated by TIKZERO+ and AUTOMATIKZ<sub>v2</sub> (LLM), demonstrating how TIKZERO+ utilizes advanced TikZ features, whereas AUTOMATIKZ<sub>v2</sub> (LLM) employs only basic, simple commands.

---

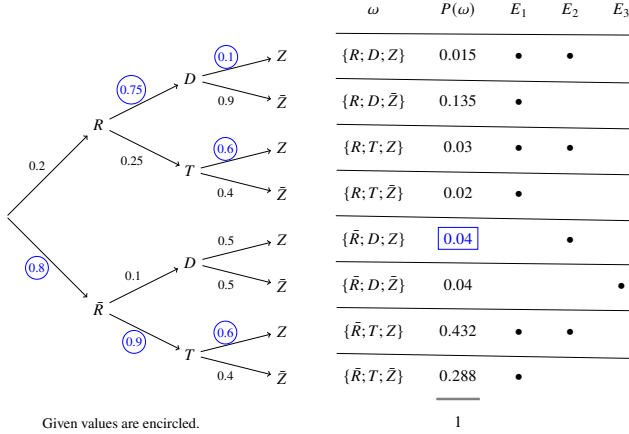
<sup>5</sup>sourced from <https://github.com/PetarV-/TikZ>, <https://github.com/janosh/tikz>, <https://tikz.net>, and <https://arxiv.org>



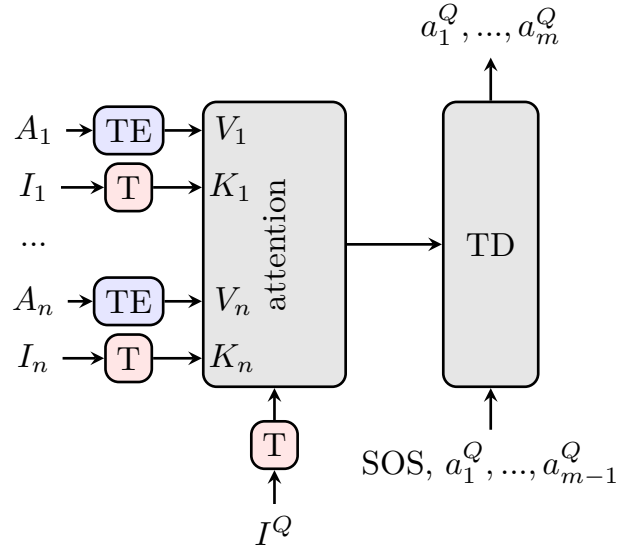
(a) A diagram representing a recurrent neural network consisting of several LSTM blocks, processing the input sequence simultaneously forwards and backwards (to exploit both directions of temporal dependence). Contains some rather tight manoeuvring.



(b) A plot comparing the distribution functions of Bose-Einstein, Boltzmann, and Fermi-Dirac statistics as a function of the reduced chemical potential  $\beta(\epsilon - \mu)$ . This visualisation highlights the differences between the three types of distribution functions, which are used to describe the behavior of particles in different statistical systems.



(c) Tree with aligned matrix. A probability tree with an aligned matrix listing the possible outcomes, their probabilities and three columns for events described in later tasks. It uses the grahdrawing library and requires LuaLaTeX.



(d) Our approach is a modified version of **meta-seq2seq**. A transformer decoder (TD) is trained to produce a sequence of actions  $a_1^Q, \dots, a_m^Q$  given a query instruction  $I^Q$ . The context are demonstrations  $(I_k, A_k)$  produced by our generative model. We use a transformer encoder-decoder (T) to encode instructions and state  $S$  and a transformer encoder (TE) to encode actions. The transformers that process instructions (pink blocks) receive state  $S$  as the input of the encoder.

Figure 5. Representative examples from  $\text{DATIKZ}_{v3}$  (also present in  $\text{DATIKZ}$  and  $\text{DATIKZ}_{v2}$ ), with permissive licenses.

Models	DSIM $\uparrow$	KID $\downarrow$	CLIP $\uparrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$	AVG $\uparrow$
TikZERO (Cos)	<b>52.829</b>	<b>5.103</b>	10.051	<b>1.603</b>	65.51	82.291	<b>64.309</b>
TikZERO (Base)	<u>52.373</u>	<u>5.225</u>	9.428	1.589	65.286	83.128	<u>63.129</u>
QWEN <sub>2.5</sub> CODER (14B)	48.352	12.988	19.761	0.229	<b>60.304</b>	<b>93.285</b>	58.894
DEEPSEEK-R1 QWEN (32B)	47.573	8.887	<u>21.201</u>	1.388	64.928	66.225	57.252
DEEPSEEK-R1 QWEN (14B)	44.616	15.43	<b>21.695</b>	0.842	<u>63.323</u>	36.11	31.102

Table 5. System-level scores  $\times 100$  for TikZERO (Cos) and additional baselines. Overall, TikZERO achieves the strongest average performance across metrics.

Models	Reference Figures					Synthetic Sketches				
	DSIM $\uparrow$	KID $\downarrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$	DSIM $\uparrow$	KID $\downarrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$
DETIKZIFY <sub>DS</sub> (7b)	75.46	0.842	2.953	56.851	84.019	67.379	0.766	1.541	59.589	84.401
DETIKZIFY <sub>v2</sub>	<b>80.503</b>	<b>0.626</b>	<b>6.105</b>	<b>54.946</b>	<b>93.326</b>	<b>74.584</b>	<b>0.751</b>	<b>3.356</b>	<b>58.32</b>	<b>93.858</b>

Table 6. System-level scores  $\times 100$  for DETIKZIFY<sub>v2</sub> and DETIKZIFY<sub>DS</sub> (7b) on both reference figures and synthetic sketches generated with ULTRASKEETCH from the test split of DATIKZ<sub>v3</sub>. Best scores are in bold, and arrows indicate metric directionality. Note that we compute DREAMSIM using updated models [68], whereas Belouadi et al. [3] used the original models in their work [67].

Data	Intv.	DSIM $\uparrow$	KID $\downarrow$	CLIP $\uparrow$	cBLEU $\uparrow$	TED $\downarrow$	MTE $\uparrow$	AVG $\uparrow$
100%	1	<b>52.771</b>	<u>5.127</u>	<u>9.949</u>	<b>1.607</b>	65.516	82.292	<b>92.411</b>
100%	2	<u>52.311</u>	5.2	<b>9.955</b>	1.484	65.473	82.588	<u>87.557</u>
100%	4	51.794	5.688	8.886	1.429	<u>65.399</u>	<b>83.988</b>	82.254
100%	8	51.59	5.933	<u>9.818</u>	1.371	65.608	<u>83.679</u>	76.545
50%	1	52.106	5.835	8.527	1.454	65.605	83.599	77.478
50%	2	52.143	<b>5.103</b>	9.315	1.393	<b>65.355</b>	82.924	85.249
50%	4	50.492	6.689	8.852	1.459	65.951	78.456	47.381
50%	8	50.093	6.738	7.999	1.379	65.963	78.923	40.816
25%	1	51.55	6.055	9.12	1.472	66.237	77.961	49.967
25%	2	51.231	6.152	8.943	1.43	65.714	77.566	54.942
25%	4	49.859	7.715	7.316	1.41	66.128	79.704	32.12
25%	8	49.179	7.764	6.495	1.434	66.009	79.9	29.774
12.5%	1	50.485	6.25	7.568	<u>1.509</u>	65.8	80.816	56.055
12.5%	2	50.152	7.129	6.353	1.275	66.045	81.05	33.817
12.5%	4	49.667	7.031	6.474	1.221	65.892	82.634	37.914
12.5%	8	48.827	8.154	5.054	1.11	65.813	80.738	16.25

Table 7. System-level scores  $\times 100$  TIKZERO (Cos) trained on varying fractions of data and intervals of cross-attention layers. Bold and underlined values denote the best and second-best scores for the whole table, respectively. Cell shading illustrates score magnitudes. Arrows indicate metric directionality.

Reference	AUTOMATikZ <sub>v2</sub>	TikZERO	TikZERO+	GPT-4o
<p>An illustration of the reduction from densest <math>k</math>-subgraph to U-RCP. On the left there is a simple undirected graph <math>G</math> with a single edge. The 2-reduced directed graph of <math>G</math> is on the right. Each vertex of <math>G</math> is replaced by <math>2 \cdot 2 = 4</math> copies with a bidirectional edge connecting any two copies of the same vertex, and an outgoing edge from each copy to the single edge-vertex <math>e</math>.</p>				
<p>Domain of dependence: The closed trapezoidal region <math>\Omega_T</math> used in the proof of Proposition 1, shaded with darker yellow; the domain of dependence of point <math>(0, \ell)</math> is <math>\Omega_\ell</math> along with the top region shaded with lighter yellow.</p>				
<p>A sketch of iterating <math>f(x) = x - 1/x</math>. Points bigger than 1 get sent to point in <math>[0, 1]</math> which then moves to a point <math>\leq -1</math> which moves to <math>[-1, 0]</math> which is sent to a point <math>\geq 1</math> and so on.</p>				
<p>A multiline diagram <math>D</math> with <math>n = 6</math> columns and <math>s = 5</math> rows, with content <math>\lambda = (5, 4, 3, 1, 0, 0)</math> and bottom row <math>\rho^{(1)}(D) = (4, 0, 1, 5, 3, 0) \in S_\lambda</math>. It has weight <math>w(D) = w_X(D) w_I(D) = x_1^3 x_3^2 x_4^4 x_5^2 x_6^2</math>.</p>				

Table 8. Alternating rows display randomly selected tuples from the caption and image similarity human evaluation task (cf. Sec. 5.3). The frames of highest and lowest rated instances are highlighted in green and red, respectively.

```

1 \documentclass[border=10pt]{standalone}
2 \usepackage{pgfplots}
3 \pgfplotsset{width=7cm,compat=1.8}
4 \begin{document}
5 \begin{tikzpicture}
6     \begin{axis}[grid=minor,view={210}{30}]
7         \addplot3+[domain=-5:5,surf]
8             {exp(-x^2-y^2)*x};
9         % connect classical points
10        \addplot3+[mark=*,only marks,blue,samples
11            at={-4.473,-1.463,1.463,4.473}]
12            {exp(-x^2-y^2)*x};
13    \end{axis}
14 \end{tikzpicture}
15 \end{document}

```

TikZERO+

```

1 \documentclass[12pt]{article}
2 \usepackage{amsmath,amsthm,amssymb,amsfonts}
3 \usepackage{tikz}
4 \usepackage{pgfplots}
5 \usetikzlibrary{arrows.meta}
6
7 \begin{document}
8
9 \begin{tikzpicture}
10 \begin{axis}[view={60}{30}]
11     \addplot3 [surf] {3*x^2 + y^2};
12 \end{axis}
13 \end{tikzpicture}
14
15 \end{document}

```

TikZERO+

```

1 \documentclass[11pt,a4paper]{article}
2 \usepackage{amsmath}
3 \usepackage{amssymb}
4 \usepackage{color}
5 \usepackage{tikz}
6
7 \begin{document}
8
9 \begin{tikzpicture}[scale=0.5]
10 \draw [thick,->] (-2,0) -- (2,0);
11 \draw [thick,->] (0,-2) -- (0,2);
12 \draw [thick] (1.8,0) arc (0:180:1.8);
13 \node [right] at (2,0) {\mathrm{Re}(\omega)};
14 \node [above] at (0,2) {\mathrm{Im}(\omega)};
15 \node [right] at (1.8,1) {\Gamma};
16 \end{tikzpicture}
17
18 \end{document}

```

AUTOMATikZ<sub>v2</sub>

```

1 \documentclass{article}
2 \usepackage{amssymb}
3 \usepackage{amsmath}
4 \usepackage{pgfplots}
5 \pgfplotsset{compat=1.14}
6 \usepackage{tikz}
7 \usetikzlibrary{arrows}
8
9 \begin{document}
10
11 \begin{tikzpicture}[scale=3]
12 \draw[fill=green!15] (1,0) -- (2,0) -- (2,2) --
13     (1,2) -- (1,0);
14 \draw[thick] (1,0) -- (1,2);
15 \draw[thick] (0,1) -- (2,1);
16 \node[below] at (1,0) {\$0\$};
17 \node[below] at (2,0) {\$1\$};
18 \node[left] at (1,2) {\$1\$};
19 \node[left] at (0,1) {\$0\$};
20 \node[above] at (1.5,1.5) {\gamma};
21 \draw[->,thick] (0,0) -- (1,0);
22 \draw[->,thick] (0,0) -- (0,1);
23 \draw[->,thick] (0,0) -- (0.5,0.5);
24 \draw[->,thick] (0,0) -- (1.5,0.5);
25 \draw[->,thick] (0,0) -- (0.5,1.5);
26 \draw[->,thick] (0,0) -- (1.5,1.5);
27 \draw[->,thick] (0,0) -- (2,0);
28 \draw[->,thick] (0,0) -- (0,2);
29 \draw[->,thick] (0,0) -- (1,1);
30 \end{tikzpicture}
31
32 \end{document}

```

AUTOMATikZ<sub>v2</sub>

Figure 6. TikZ programs generated by TikZERO+ (top) and AUTOMATikZ<sub>v2</sub> (LLM; bottom) corresponding to the figures shown in the first row of Fig. 1 in the same order.