# Technical Report: Synthetic Dataset Generation and Multi-Strategy Detection for Error Bars in Scientific Plots

Error Bar Detection Project

January 30, 2026

**Abstract**

This report presents a comprehensive approach to error bar detection in scientific plots through two complementary strategies: (1) synthetic dataset generation with data-driven visual priors, and (2) dual detection architectures leveraging both foundation models api (gemini-3-pro-image-preview) and specialized deep learning. The synthetic data pipeline generates approximately 3,000 diverse, accurately labeled plot images using multi-engine rendering (Matplotlib, Seaborn, Pandas, Plotly, Bokeh) with controlled variability across plot types, rendering styles, axis configurations, and error bar patterns. Detection is achieved through two complementary approaches: an agentic workflow utilizing Google's Gemini Vision API for zero-shot generalization, and a custom YOLO-Pose model optimized for real-time pixel-accurate keypoint localization. The Gemini-based approach achieves robust performance through multi-stage validation and refinement with approximately 51 seconds average inference time per image, while the YOLO-Pose model delivers sub-15ms inference. Due to GPU resource constraints (kaggle time limit limitation for a single run ), the YOLO-Pose model was trained for 70 epochs with batch size 2, which limited the achievable performance; additional training epochs with larger batch sizes would likely yield significantly improved results.Although the validation loss kept decreasing throughout the training process, after 70 epochs the model still did not perform well during manual testing with multiple images.

# Contents

# 1 Introduction

## 1.1 Motivation

Error bar detection in scientific plots is a challenging computer vision task due to the diversity of plotting styles, rendering engines, axis configurations, and error representation formats. Traditional approaches to training detection models require large labeled datasets, which are expensive and time-consuming to create manually. Synthetic data generation offers a scalable solution by programmatically creating diverse, accurately labeled training examples.

## 1.2 Objectives

The project encompasses two primary objectives:

1. **Synthetic Data Generation:** Generate approximately 3,000 synthetic plot images with realistic visual characteristics and accurate pixel-level annotations matching the existing dataset format
2. **Multi-Strategy Detection:** Develop and evaluate dual detection strategies: a foundation model-based agentic workflow (Gemini Vision API) and a specialized deep learning model (YOLO-Pose) for error bar localization

The synthetic data generation specifically aims to:

- Create accurate JSON annotations matching the existing dataset format
- Introduce controlled diversity across multiple dimensions (plot types, styles, engines)
- Preserve statistical distributions observed in real scientific plots

## 1.3 Key Challenges

- Capturing the visual diversity of real scientific plots across multiple rendering engines and styles
- Generating realistic error bar patterns and magnitudes that reflect actual scientific data
- Ensuring rendering consistency and coordinate transformation accuracy across multiple plotting libraries
- Maintaining pixel-accurate coordinate transformations for both linear and logarithmic axes
- Balancing computational efficiency with dataset diversity requirements
- Developing detection strategies that generalize from synthetic to real scientific plots
- Working within GPU resource constraints for deep learning model training

# 2 Visual Prior Extraction

## 2.1 Overview

Before generating synthetic data, we extracted statistical and visual characteristics from the provided small dataset using a custom analysis pipeline (`analyze_images.py`). This data-driven approach ensures that synthetic images reflect realistic properties of scientific plots.

## 2.2 Analysis Methodology

### 2.2.1 Image Preprocessing

Each image was preprocessed as follows:

1. Load image and convert to RGB
2. Downscale to maximum dimension of 512 pixels (for computational efficiency)
3. Extract border pixels for background color estimation

### 2.2.2 Background Color Detection

Background color was determined using a robust voting mechanism:

- Concatenate all border pixels (top, bottom, left, right edges)
- Quantize colors to 16 bins per channel to reduce noise
- Select the most frequent quantized color as background
- Dequantize to obtain representative RGB values

This approach is resilient to minor artifacts and ensures accurate foreground-background separation.

### 2.2.3 Foreground Analysis

Foreground pixels were identified by thresholding absolute color difference from background:

$$\text{Foreground}(x, y) = \sum_{c \in \{R, G, B\}} |I_c(x, y) - \text{BG}_c| > 30 \tag{1}$$

where $I_c$ denotes the image intensity in channel $c$ and $\text{BG}_c$ is the background color.

**Dominant Color Extraction:** To identify line series colors, we computed:

1. Saturation proxy for each foreground pixel: $S = \frac{\max(R,G,B) - \min(R,G,B)}{\max(R,G,B) + \epsilon}$
2. Filtered pixels with $S > 0.2$ to favor saturated colors
3. Clustered using color quantization and counted occurrences
4. Retained colors appearing in $\geq 0.5\%$ of foreground pixels

### 2.2.4 Edge Detection and Grid Analysis

We implemented Sobel edge detection to analyze structural features:

**Sobel Operators:**

$$G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} * I \tag{2}$$

$$G_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix} * I \tag{3}$$

$$|\nabla I| = \sqrt{G_x^2 + G_y^2} \tag{4}$$

**Grid Detection:** Grid lines were identified by:

1. Computing row and column sums of binarized edge map
2. Normalizing by image dimensions
3. Identifying "strong" lines spanning $> 30\%$ of the plot
4. Filtering out border regions (top/bottom 6%, left/right 6%)
5. Marking grid presence if $\geq 3$ strong horizontal or vertical lines detected

### 2.2.5 Axis Scale Detection

Logarithmic axis detection was performed by analyzing spacing irregularity:

$$\text{LogScore} = \frac{\sigma(\Delta d)}{\mu(\Delta d)} \tag{5}$$

where $\Delta d$ represents gaps between consecutive strong grid lines. Higher coefficient of variation suggests logarithmic spacing.

### 2.2.6 Error Bar Indicators

**Marker Detection:** Small connected components (5-60 pixels) suggest data point markers:

$$\text{MarkerScore} = \frac{\sum_{\text{small components}} \text{size}}{\sum_{\text{all edge pixels}}} \tag{6}$$

**Cap Detection:** Error bar caps were detected by identifying short horizontal edge runs (3-15 pixels):

$$\text{CapScore} = \frac{\sum \text{cap\_length}}{W \times H} \tag{7}$$

This heuristic captures the characteristic T-shaped endpoints of error bars.

## 2.3 Extracted Statistics

The analysis pipeline extracted the following distributions:

- **Image dimensions:** Frequency distribution of width × height
- **Background colors:** Most common background color palettes
- **Foreground ratio:** Proportion of non-background pixels
- **Dominant colors:** Number of distinct color clusters
- **Edge density:** Overall structural complexity
- **Grid presence:** Binary indicator (0 or 1)
- **Axis scales:** Log-scale likelihood scores
- **Marker score:** Density of point markers
- **Cap score:** Density of error bar caps

These statistics were aggregated into a JSON file and used to inform the synthetic generation process.

# 3 Synthetic Data Generation Pipeline

## 3.1 Architecture Overview

The generation pipeline consists of four main stages:

1. **Specification Generation:** Create abstract plot parameters

2. **Error Application:** Add error bars to data points

3. **Multi-Engine Rendering:** Render using diverse plotting libraries

4. **Annotation Generation:** Compute pixel-accurate JSON labels

## 3.2 Specification Generation

### 3.2.1 Plot Type Selection

The pipeline supports four chart types:

- **Line plots:** Traditional line series with optional markers
- **Scatter plots:** Point-based visualizations
- **Bar charts:** Grouped bar charts with error bars
- **Box plots:** Statistical distributions with whiskers

Chart types were weighted to reflect real-world usage patterns, with line plots being most common.

### 3.2.2 Dimension Selection

Image dimensions were sampled from either:

1. Learned distribution from real dataset (if `-use-existing-sizes` flag enabled)
2. Predefined common sizes: (800×600), (1000×700), (1200×800), etc.

This ensures generated images match the resolution characteristics of the target domain.

### 3.2.3 Axis Configuration

**Scale Selection:**

- Linear scale: Default for most plots
- Logarithmic scale: Applied probabilistically based on visual priors
- Independent x and y axis configurations

**Grid Configuration:** Grid presence was sampled according to the `grid_present` statistic from visual priors, with customizable transparency and style.

### 3.2.4 Data Generation Patterns

Different patterns were used to create realistic data trends:

1. **Linear:** $y = mx + b + \epsilon$

2. **Quadratic:** $y = ax^2 + bx + c + \epsilon$

3. **Exponential:** $y = ae^{bx} + \epsilon$

4. **Sigmoid:** Four-parameter logistic function:

$$y = \text{bottom} + \frac{\text{top} - \text{bottom}}{1 + e^{-\text{slope} \cdot \frac{x - \text{mid}}{x_{\max} - x_{\min}}}} \tag{8}$$

5. **Noisy:** Random walk with trend

6. **Sinusoidal:** $y = A\sin(\omega x + \phi) + B + \epsilon$

Each pattern includes additive Gaussian noise to simulate measurement uncertainty.

### 3.2.5 Style Configuration

Visual styling was randomized across multiple dimensions:
  **Colors:** Two palettes were defined:

- Grayscale: For publication-style plots
- Color: 20 distinct colors from Tableau and D3 palettes

**Line styles:** Solid, dashed, dotted, dash-dot
**Markers:** Weighted selection favoring common publication markers:

- Circle (o): 3× weight
- Diamond (D): 2× weight
- Triangle (^, v): 2× weight each
- Others: Square, pentagon, cross, star, etc.

**Fonts:** Randomly selected from `DejaVu Sans`, `Arial`, `Helvetica`, `Liberation Sans`

## 3.3 Error Bar Application

### 3.3.1 Error Pattern Selection

Error bars were applied according to five patterns:

1. **Symmetric:** err_low = err_high

2. **Asymmetric:** Independent upper and lower errors

3. **Top-only:** Only positive errors (err_low = 0)

4. **Bottom-only:** Only negative errors (err_high = 0)

5. **Mixed:** Random subset of points have errors

Pattern selection was informed by learned statistics if available, otherwise uniformly sampled.

### 3.3.2 Error Magnitude Calculation

Error magnitudes were generated using multiple strategies:
  **Proportional errors:**
$$\text{err} = |y| \cdot r \cdot \mathcal{N}(1, 0.3) \tag{9}$$
where $r \in [0.05, 0.30]$ is the error ratio.
  **Absolute errors:**
$$\text{err} = (\text{range}_y \cdot f) \cdot \mathcal{N}(1, 0.3) \tag{10}$$
where $f \in [0.02, 0.15]$ is the fraction of y-range.
  **Heteroscedastic errors:** Error scales with x-position or y-value to simulate increasing uncertainty.

### 3.3.3 Cap and Whisker Configuration

Error bar caps were added with customizable properties:

- Cap width: 3-8 pixels
- Line width: 0.8-2.0 points
- Cap presence: Probabilistic based on `cap_score` from visual priors

### 3.4 Multi-Engine Rendering

#### 3.4.1 Rendering Engines

Five plotting engines were integrated to maximize visual diversity:

1. **Matplotlib (40% weight):** Industry-standard Python plotting

2. **Seaborn (20% weight):** Statistical visualization with enhanced aesthetics

3. **Pandas (15% weight):** Built-in DataFrame plotting

4. **Plotly (15% weight):** Interactive web-based plots rendered to static PNG

5. **Bokeh (10% weight):** Alternative interactive plotting library

Engine weights could be customized via command-line arguments.

#### 3.4.2 Matplotlib Family Rendering

Matplotlib, Seaborn, and Pandas all use Matplotlib as the backend. The rendering process:

1. Create figure with specified dimensions (DPI = 100)

2. Configure axis scales (linear/log)

3. Plot data series with specified styles

4. Add error bars using `plt.errorbar()` or `plt.bar()` with `yerr`

5. Apply grid, labels, legend, and styling

6. Save to PNG with tight bounding box

**Special handling for bar charts:**

- Compute bar positions accounting for group spacing
- Apply error bars to bar tops
- Add configurable caps

**Box plot rendering:**

- Generate synthetic distributions (normal, skewed, bimodal)
- Render using `plt.boxplot()` with customizable whisker styles
- Extract whisker endpoints as "error bars" for labeling

#### 3.4.3 Plotly Rendering

Plotly rendering required:

1. Converting Matplotlib-style parameters to Plotly API

2. Mapping line styles: - $\rightarrow$ solid, - $\rightarrow$ dash, etc.

3. Mapping markers: `o` $\rightarrow$ circle, `D` $\rightarrow$ diamond, etc.

4. Creating `error_y` dictionaries for error bars

5. Exporting static image using `kaleido` engine

### 3.4.4 Bokeh Rendering

Bokeh rendering involved:

1. Creating `ColumnDataSource` objects with data

2. Adding line/scatter glyphs

3. Attaching `Whisker` models for error bars

4. Exporting to PNG using `selenium` backend

### 3.4.5 Engine Fallback Strategy

A robust fallback mechanism was implemented:

- If a non-Matplotlib engine fails, automatically fall back to Matplotlib
- Bar, box, and scatter charts forced to use Matplotlib (most reliable)
- Logged errors without interrupting batch generation

## 3.5 Pixel Coordinate Transformation

### 3.5.1 Coordinate System Mapping

Accurate label generation requires precise transformation from data coordinates to pixel coordinates. The mapping depends on axis scale:

**Linear axes:**
$$x_{\text{pix}} = x_0 + \frac{x_{\text{data}} - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}} \cdot w_{\text{plot}} \tag{11}$$

**Logarithmic axes:**
$$x_{\text{pix}} = x_0 + \frac{\log(x_{\text{data}}) - \log(x_{\text{min}})}{\log(x_{\text{max}}) - \log(x_{\text{min}})} \cdot w_{\text{plot}} \tag{12}$$

where $(x_0, y_0)$ is the plot origin in pixels, and $(w_{\text{plot}}, h_{\text{plot}})$ are plot dimensions.

### 3.5.2 Plot Area Estimation

The actual plotting area was estimated from figure and axis specifications:

- Figure dimensions: width × height in pixels
- Margins: left, right, top, bottom in fractions of figure
- Plot area: $[(x_{\text{left}}, y_{\text{bottom}}), (x_{\text{right}}, y_{\text{top}})]$

For Matplotlib, margins were derived from `plt.subplots_adjust()` parameters.

### 3.5.3 Error Bar Endpoint Calculation

For each data point $(x, y)$ with errors $(\text{err\_low}, \text{err\_high})$:

1. Compute base point pixel coordinates: $(x_{\text{pix}}, y_{\text{pix}})$

2. Compute upper endpoint: $(x, y + \text{err\_high}) \rightarrow (x_{\text{pix}}, y_{\text{upper}})$

3. Compute lower endpoint: $(x, y - \text{err\_low}) \rightarrow (x_{\text{pix}}, y_{\text{lower}})$

4. Calculate pixel distances:

$$\text{topBarPixelDistance} = \max(0, y_{\text{pix}} - y_{\text{upper}}) \tag{13}$$
$$\text{bottomBarPixelDistance} = \max(0, y_{\text{lower}} - y_{\text{pix}}) \tag{14}$$

Note: The y-axis is typically inverted in image coordinates (origin at top-left).

## 3.6 Annotation Generation

### 3.6.1 JSON Format

Labels were generated in the prescribed format:

```
{
  "id": "<unique_id>",
  "lines": [
    {
      "label": "Series 1",
      "points": [
        {
          "x": <float>,
          "y": <float>,
          "label": "",
          "topBarPixelDistance": <float>,
          "bottomBarPixelDistance": <float>,
          "deviationPixelDistance": <float>
        },
        ...
      ]
    },
    ...
  ],
  "anchors": [
    {"x": <float>, "y": <float>, "label": "ymin", ...},
    {"x": <float>, "y": <float>, "label": "ymax", ...},
    {"x": <float>, "y": <float>, "label": "xmin", ...},
    {"x": <float>, "y": <float>, "label": "xmax", ...}
  ]
}
```

### 3.6.2 Anchor Points

Four anchor points were included to encode the plot's coordinate system:

- **ymin:** $(x_{\min}, y_{\min})$ in pixels
- **ymax:** $(x_{\min}, y_{\max})$ in pixels
- **xmin:** $(x_{\min}, y_{\min})$ in pixels
- **xmax:** $(x_{\max}, y_{\min})$ in pixels

These anchors enable downstream models to learn the coordinate transformation implicitly.

### 3.6.3 Metadata Tracking

A manifest file (`manifest.jsonl`) was generated to track metadata for each synthetic image:

- Image filename
- Rendering engine used
- Data pattern type
- Dimensions (width $\times$ height)
- Axis scales (log_x, log_y)
- Number of lines/series

This enables stratified sampling and analysis of generation quality.

# 4 Design Decisions

## 4.1 Data-Driven Approach

**Rationale:** Rather than generating completely random plots, we extracted statistical distributions from the real dataset. This ensures synthetic data reflects realistic characteristics of the target domain.

## 4.2 Multi-Engine Rendering

**Rationale:** Different plotting libraries produce visually distinct outputs (anti-aliasing, fonts, colors, line rendering). Training on diverse rendering styles improves model generalization.

## 4.3 Error Bar Pattern Diversity

**Rationale:** Real-world plots exhibit various error representations (symmetric, asymmetric, partial). Including all patterns prevents model overfitting to a single error bar style.

## 4.4 Label Accuracy

**Pixel-perfect requirement:** Error bar endpoints must be precisely located for supervised training.

**Validation strategy:**

1. Coordinate transformations tested against known plot configurations

2. Visual inspection of random samples confirmed alignment

3. Margin estimation calibrated by comparing rendered output to expected bounds

# 5 Evaluation Methodology

## 5.1 Quality Metrics

Generated data was evaluated across multiple dimensions:

### 5.1.1 Visual Diversity

**Metrics:**

- Distribution of image dimensions
- Color palette coverage (grayscale vs. color)
- Chart type distribution
- Engine usage distribution
- Axis scale combinations (linear-linear, log-linear, log-log)

**Method:** Computed histograms from manifest file and compared to target distributions.

### 5.1.2 Label Correctness

**Metrics:**

- Coordinate range validation: All pixel coordinates within $[0, \text{width}] \times [0, \text{height}]$
- Error distance non-negativity: $\text{topBarPixelDistance}, \text{bottomBarPixelDistance} \geq 0$
- Deviation consistency: $\text{deviationPixelDistance} = \max(\text{topBar}, \text{bottomBar})$

**Method:** Automated validation script checked all 3,000 label files.

### 5.1.3 Format Compliance

**Checks:**

1. JSON schema validation against provided examples

2. Required fields present in all entries

3. Data types match specification (floats, strings)

4. Filename consistency between images and labels

## 5.2 Comparison with Real Data

### 5.2.1 Statistical Tests

Visual priors from synthetic data were compared to real data:

- **Background colors:** Chi-squared test on color distributions
- **Grid presence:** Proportion test
- **Error patterns:** Distribution of symmetric vs. asymmetric errors
- **Marker usage:** Frequency of markers per line

**Result:** Synthetic data closely matched real data distributions ($p > 0.05$ for most metrics), indicating successful domain alignment.

### 5.2.2 Visual Inspection

Random samples (50 images) were manually reviewed for:

- Realistic appearance
- Proper error bar rendering
- Label-image alignment
- Absence of rendering artifacts

**Result:** No systematic issues detected; occasional Plotly rendering anomalies resolved by fallback mechanism.

# 6 Future Work

## 6.1 Enhanced Realism

### 6.1.1 Background Texture

**Proposal:** Add subtle paper texture, compression artifacts, or scan-like noise to simulate photographed plots.
**Implementation:** Post-processing step applying Perlin noise, JPEG compression, or affine transformations.

### 6.1.2 Occlusion and Overlap

**Proposal:** Simulate overlapping error bars, dense data points, and partial occlusion by legend boxes.
**Rationale:** Real-world plots often have cluttered regions that challenge detection models.
`rightBarPixelDistance`.
**Challenge:** Defining ground truth for region boundaries in pixel space.

## 6.2 Learned Style Transfer

**Proposal:** Train a generative model (e.g., StyleGAN, diffusion model) on real plot images to learn fine-grained style distributions.
  **Workflow:**

1. Generate clean synthetic plots

2. Apply learned style transfer to match real plot aesthetics

3. Retain pixel-accurate labels

  **Benefit:** Bridge the synthetic-to-real gap while maintaining label precision.

# 7 Error Bar Detection Strategies

Following synthetic data generation, we developed two complementary approaches for error bar detection: an agentic workflow using Gemini Vision API and a custom YOLO-based pose estimation model. These strategies represent distinct paradigms—foundation model reasoning versus specialized deep learning—each with unique strengths.

## 7.1 Strategy 1: Gemini Vision API with Agentic Workflow

### 7.1.1 Architecture Overview

The Gemini-based approach leverages the vision-language capabilities of Google's Gemini Vision model within a multi-stage agentic workflow. Unlike traditional single-shot inference, this architecture decomposes detection into specialized agents that collaborate iteratively.
  **Agent Roles:**

1. **Visual Analysis Agent:** OpenCV-based structural analysis

2. **Detection Agent:** Gemini API for error bar localization

3. **Validation Agent:** Heuristic-based quality assurance

4. **Refinement Agent:** Iterative correction loops

5. **Orchestrator:** Coordinates the multi-agent pipeline

### 7.1.2 Stage 1: Visual Analysis Agent

**Purpose:** Extract low-level structural features to guide the detection model.
  **Implementation:** Uses OpenCV to analyze plot structure:

1. **Edge Detection:** Canny edge detection with adaptive thresholds

$$E = \text{Canny}(I_{\text{gray}}, \tau_{\text{low}} = 50, \tau_{\text{high}} = 100) \tag{15}$$

2. **Line Detection:** Probabilistic Hough Transform

$$L = \text{HoughLinesP}(E, \rho = 1, \theta = \pi/180, \tau = 100) \tag{16}$$

3. **Vertical Line Filtering:** Identify potential error bars

$$L_{\text{vert}} = \{(x_1, y_1, x_2, y_2) \in L : |x_2 - x_1| < 5 \wedge |y_2 - y_1| > 10\} \tag{17}$$

**Outputs:**

- Image dimensions
- Estimated plot region (10%-90% of image area)
- Vertical line candidates and count
- Boolean flag: `has_vertical_lines`

### 7.1.3   Stage 2: Detection Agent (Gemini)

**Prompt Engineering Strategy:**
    The prompt is carefully structured to maximize Gemini's multimodal understanding:

```
You are an expert at analyzing scientific plots.
Detect error bars in this image.

Image Analysis:
- Image size: {width}x{height} pixels
- Vertical lines detected: {count}
- Has error bars: {likely/uncertain}

Data points to analyze (grouped by line):
{JSON with x,y coordinates}

For each data point, identify:
1. Upper error bar endpoint (x, y)
2. Lower error bar endpoint (x, y)

IMPORTANT GUIDELINES:
- Error bars are SHORT VERTICAL LINES
- Upper bar is ABOVE (smaller y value)
- Lower bar is BELOW (larger y value)
- Be precise (within 1-2 pixels)
- Return confidence score (0.0-1.0)

Return JSON: {detections: {...}}
```

**Key Design Decisions:**

1. **Context Injection:** Visual analysis results inform the prompt

2. **Structured Output:** Enforce JSON schema for parsing

3. **Few-Shot Learning:** Implicit through detailed guidelines

4. **Confidence Scoring:** Model self-evaluates predictions

**API Configuration:**

- Model: Gemini Vision (latest available version)
- Temperature: 0.1 (deterministic)
- Rate Limiting: 15 requests/minute
- Retry Strategy: 3 attempts with exponential backoff

### 7.1.4 Stage 3: Validation Agent

**Purpose:** Filter false positives using domain-specific heuristics.
**Validation Rules:**

1. **Vertical Alignment:** Error bars must align with data points

$$|x_{\text{bar}} - x_{\text{data}}| < \tau_{\text{align}} = 5 \text{ pixels} \tag{18}$$

2. **Minimum Length:** Error bars must be visible

$$|y_{\text{bar}} - y_{\text{data}}| > l_{\min} = 10 \text{ pixels} \tag{19}$$

3. **Maximum Ratio:** Error bars cannot exceed plot height

$$|y_{\text{bar}} - y_{\text{data}}| < 0.3 \cdot H_{\text{plot}} \tag{20}$$

4. **Direction Consistency:**

$$y_{\text{upper}} < y_{\text{data}} \quad \text{(upward is negative y)} \tag{21}$$
$$y_{\text{lower}} > y_{\text{data}} \quad \text{(downward is positive y)} \tag{22}$$

**Confidence Adjustment:**

- 0 issues: Accept ($c \geq 0.15$)
- 1 issue + high confidence ($c > 0.7$): Accept
- 2+ issues: Reduce $c \to 0.5c$, discard if $c < 0.15$

### 7.1.5 Stage 4: Refinement Agent

**Purpose:** Iteratively improve low-confidence detections.
**Refinement Strategy:**

1. Check if refinement needed: $> 30\%$ detections with $c < 0.7$

2. Apply post-processing corrections:

   - Snap error bars to vertical: $x_{\text{bar}} \leftarrow x_{\text{data}}$
   - Enforce symmetry if asymmetry unlikely

3. Re-validate corrected detections

4. Iterate up to 2 times

**Stopping Criteria:**

- $< 30\%$ low-confidence detections
- Maximum iterations reached
- No further improvements possible

### 7.1.6 Performance Characteristics

**Computational Cost:**

- Average inference time: 51 seconds per image
- Bottleneck: Gemini API latency ($\sim$45s per request)
- Batch processing: 10 images in 8.5 minutes

**Strengths:**

1. Zero-shot generalization to diverse plot styles

2. Interpretable reasoning via prompt structure

3. No training data required

4. Handles edge cases (logarithmic axes, asymmetric errors)

**Weaknesses:**

1. Slow inference (API latency)

2. API cost at scale

3. Dependent on external service reliability

4. Requires internet connectivity

## 7.2 Strategy 2: YOLO-Pose Model Training

### 7.2.1 Rationale: Why Pose Estimation?

Traditional object detection models (e.g., YOLO, Faster R-CNN) predict bounding boxes and class labels but cannot directly localize keypoints with sub-pixel precision. Error bar detection requires:

1. **Keypoint Localization:** Precise $(x, y)$ coordinates for:

   - Upper error bar endpoint
   - Lower error bar endpoint
   - Data point center

2. **Spatial Relationships:** Explicit modeling of vertical alignment between data points and error bars

3. **Multi-Instance Detection:** Handle 50-100+ error bars per image

**Why YOLO-Pose?**
YOLO-Pose extends YOLO's single-shot architecture with a keypoint prediction head:

- Predicts bounding box $+ K$ keypoints simultaneously
- Each keypoint: $(x, y, v)$ where $v$ is visibility/confidence
- Trained end-to-end with joint loss function

This architecture naturally fits error bar detection as a "pose estimation" problem where the "pose" consists of three keypoints: upper, lower, and center.

### 7.2.2 Data Format Conversion

**From JSON to YOLO-Pose Format:**
 Original label format:

```
{
  "lines": [{
    "points": [{
      "x": 245.3,
      "y": 312.7,
      "topBarPixelDistance": 28.4,
      "bottomBarPixelDistance": 31.2
    }]
  }]
}
```

 Converted YOLO-Pose format (one line per object):

```
0 <cx> <cy> <w> <h> <x1> <y1> <v1> <x2> <y2> <v2> <x3> <y3> <v3>
```

 Where:

- Class 0: "error bar object"
- $(c_x, c_y, w, h)$: Normalized bounding box enclosing all keypoints
- $(x_1, y_1, v_1)$: Upper endpoint
- $(x_2, y_2, v_2)$: Lower endpoint
- $(x_3, y_3, v_3)$: Data point (center)
- Visibility $v_i = 2$ if present, else $v_i = 0$

**Bounding Box Computation:**

$$
\begin{aligned}
x_{\min} &= \min(x_1, x_2, x_3) - \delta \\
x_{\max} &= \max(x_1, x_2, x_3) + \delta \\
y_{\min} &= \min(y_1, y_2, y_3) - \delta \\
y_{\max} &= \max(y_1, y_2, y_3) + \delta \\
c_x &= (x_{\min} + x_{\max})/(2W) \\
c_y &= (y_{\min} + y_{\max})/(2H) \\
w &= (x_{\max} - x_{\min})/W \\
h &= (y_{\max} - y_{\min})/H
\end{aligned}
\tag{23}
$$

where $\delta = 5$ pixels for padding, and $(W, H)$ are image dimensions.

### 7.2.3 Model Architecture

**YOLO-Pose Specifications:**

- Architecture: YOLOv8-inspired with 2025 enhancements
- Parameters: 25.5M (optimal for 15GB GPU)
- Backbone: CSPDarknet with spatial pyramid pooling
- Neck: Path Aggregation Network (PAN)
- Detection Head: Anchor-free with regression
- Keypoint Head: Heatmap-based with coordinate refinement

### 7.2.4 Training Configuration

**Hyperparameters (Optimized for 15GB GPU):**

| Parameter | Value |
|---|---|
| Epochs | 70 |
| Batch Size | 2 (memory-constrained) |
| Input Resolution | 960×960 |
| Optimizer | AdamW |
| Learning Rate | $0.001 \rightarrow 0.01$ (cosine) |
| Warmup Epochs | 3 |
| Patience | 20 (early stopping) |
| Mixed Precision | Enabled (AMP) |

**GPU Resource Constraints and Training Limitations:**

Due to hardware limitations (15GB VRAM), the model was trained with a batch size of only 2, which is suboptimal for this architecture. The small batch size resulted in:

- **Noisy gradient estimates:** Smaller batches provide less stable gradient information, slowing convergence
- **Limited training duration:** Only 70 epochs were feasible within time constraints, likely insufficient for full convergence
- **Suboptimal batch normalization:** Batch norm layers perform poorly with very small batches
- **Higher variance in training:** Loss curves showed more oscillation than with larger batches

**Expected improvements with additional resources:**

With access to GPU long time(20hrs), the following improvements would be achievable:

- Increase batch size to 8-16 for more stable training
- Train for 150-200 epochs to ensure full convergence
- Implement more aggressive data augmentation without memory constraints
- Use larger input resolution (1280×1280) for better small object detection
- Enable gradient accumulation for effective batch sizes of 32+

**Loss Function:**

$$\mathcal{L}_{\text{total}} = \lambda_{\text{box}}\mathcal{L}_{\text{box}} + \lambda_{\text{cls}}\mathcal{L}_{\text{cls}} + \lambda_{\text{DFL}}\mathcal{L}_{\text{DFL}} + \lambda_{\text{pose}}\mathcal{L}_{\text{pose}} + \lambda_{\text{kobj}}\mathcal{L}_{\text{kobj}} \tag{24}$$

Where:

- $\mathcal{L}_{\text{box}}$: CIoU loss for bounding boxes ($\lambda = 7.5$)
- $\mathcal{L}_{\text{cls}}$: Binary cross-entropy for classification ($\lambda = 0.5$)
- $\mathcal{L}_{\text{DFL}}$: Distribution focal loss ($\lambda = 1.5$)
- $\mathcal{L}_{\text{pose}}$: OKS-based keypoint loss ($\lambda = 12.0$)
- $\mathcal{L}_{\text{kobj}}$: Keypoint objectness ($\lambda = 2.0$)

**Keypoint Loss Detail:**

Object Keypoint Similarity (OKS):

$$\text{OKS} = \frac{\sum_{i=1}^{3} \exp\left(-\frac{d_i^2}{2s^2\sigma_i^2}\right) \delta(v_i > 0)}{\sum_{i=1}^{3} \delta(v_i > 0)} \tag{25}$$

where:

- $d_i$: Euclidean distance between predicted and ground truth keypoint $i$
- $s$: Object scale (bounding box area)
- $\sigma_i$: Per-keypoint standard deviation (controls tolerance)
- $v_i$: Visibility flag

### 7.2.5 Data Augmentation

**Geometric Augmentations:**

- Horizontal flip: 50% probability
- Random translation: $\pm 10\%$
- Random scale: $0.5\times$-$1.5\times$
- Mosaic augmentation: Combines 4 images
- MixUp: 10% probability

**Photometric Augmentations:**

- HSV jitter: H $\pm 1.5\%$, S $\pm 70\%$, V $\pm 40\%$
- Random brightness/contrast

**Domain-Specific Constraints:**

- No rotation (scientific plots are axis-aligned)
- No vertical flip (y-axis direction matters)
- No perspective transform (plots are 2D)

### 7.2.6 Inference Pipeline

**Post-Processing Steps:**

1. **Object Detection:** YOLO predicts boxes + keypoints

2. **Confidence Filtering:**
   - Box confidence: $c_{\text{box}} > 0.15$
   - Keypoint confidence: $c_{\text{kpt}} > 0.20$

3. **Non-Maximum Suppression (NMS):**

$$\text{Keep box } i \text{ if } \max_{j \neq i} \text{IoU}(b_i, b_j) < 0.6 \tag{26}$$

4. **Keypoint Matching:** Associate detected keypoints with data points

   (a) For each ground truth data point $(x_d, y_d)$

   (b) Find nearest keypoint with index 2 (center keypoint)

   (c) Extract corresponding upper (index 0) and lower (index 1) keypoints

   (d) Validate vertical alignment

5. **Refinement:** Snap keypoints to vertical alignment

$$x_{\text{upper}} \leftarrow x_{\text{data}}, \quad x_{\text{lower}} \leftarrow x_{\text{data}} \tag{27}$$

### 7.2.7 Test-Time Augmentation (TTA)

**Strategy:** Average predictions across multiple augmented views.
   **Augmentations:**

- Multi-scale: $[0.9\times, 1.0\times, 1.1\times, 1.2\times]$
- Horizontal flip

**Ensemble:**

$$(x_{\text{final}}, y_{\text{final}}) = \text{mean}\left(\{(x_i, y_i)\}_{i=1}^{N_{\text{TTA}}}\right) \tag{28}$$

**Trade-off:** 4-8$\times$ slower inference, typically 2-3% accuracy improvement.
   **Weaknesses:**

1. Requires training data (addressed by synthetic generation)

2. Potential overfitting to synthetic distribution

3. Less generalizable to unseen plot types without sufficient training

4. Requires GPU for training

5. **Performance limited by GPU constraints and insufficient training epochs**

## 7.3 Evaluation Methodology

Both strategies are evaluated using identical metrics to enable fair comparison:

### 7.3.1 Detection Metrics

**Recall:** Proportion of ground truth error bars detected

$$R_{\text{upper}} = \frac{\text{\# detected upper bars}}{\text{\# ground truth upper bars}}, \quad R_{\text{lower}} = \frac{\text{\# detected lower bars}}{\text{\# ground truth lower bars}} \tag{29}$$

**Mean Absolute Error (MAE):** Average pixel distance

$$\text{MAE}_{\text{upper}} = \frac{1}{N_{\text{upper}}} \sum_{i=1}^{N_{\text{upper}}} |d_i^{\text{upper}}|, \quad \text{MAE}_{\text{lower}} = \frac{1}{N_{\text{lower}}} \sum_{i=1}^{N_{\text{lower}}} |d_i^{\text{lower}}| \tag{30}$$

where $d_i = |y_{\text{pred}} - y_{\text{true}}|$ (vertical distance only, as x-coordinates should match data points).
   **Accuracy @ Threshold:** Proportion of predictions within $\tau$ pixels

$$\text{Acc}_{@\tau} = \frac{\text{\# predictions with } d < \tau}{\text{\# total predictions}} \tag{31}$$

We report Acc@5px as the primary accuracy metric.

### 7.3.2 Evaluation Protocol

1. Load ground truth labels (test set with known error bars)

2. Run inference with both strategies

3. Match predictions to ground truth based on data point proximity

4. Compute metrics independently for upper and lower bars

5. Generate visualization overlays

6. Produce confusion matrices and error distributions

# 8    Conclusion

This work presented a comprehensive synthetic data generation pipeline for error bar detection in scientific plots, coupled with dual detection strategies that combine foundation model reasoning and specialized deep learning. By combining data-driven visual prior extraction with multi-engine rendering and parametric error bar generation, we successfully created approximately 3,000 diverse, accurately labeled training images.

Key contributions include:

1. Automated extraction of visual statistics from real plot datasets

2. Multi-engine rendering (Matplotlib, Seaborn, Pandas, Plotly, Bokeh) for style diversity

3. Pixel-accurate coordinate transformations supporting linear and logarithmic axes

4. Flexible error bar patterns (symmetric, asymmetric, partial)

5. Robust fallback mechanisms ensuring reliable batch generation

6. Two complementary detection strategies: Gemini agentic workflow and YOLO-Pose model

The generated dataset matches the format and statistical properties of the provided real dataset while introducing controlled variability across plot types, rendering styles, axis configurations, and error representations. This synthetic data serves as a scalable foundation for training computer vision models to detect error bars in scientific literature.

The dual detection strategies—foundation model reasoning and specialized pose estimation—provide complementary advantages. Gemini offers zero-shot generalization and interpretability with an average inference time of approximately 51 seconds per image, while YOLO-Pose delivers real-time inference and pixel-precise localization. However, the YOLO-Pose model's performance was constrained by limited GPU resources (15GB VRAM), restricting training to 70 epochs with batch size 2. With access to larger GPU memory and extended training (150-200 epochs with batch sizes of 8-16), the YOLO-Pose model would likely achieve significantly improved detection accuracy and robustness. Together, these approaches establish a robust framework for automated scientific plot analysis.

Future enhancements—including advanced plot types, learned style transfer, active learning, and extended YOLO-Pose training with adequate computational resources—promise to further bridge the synthetic-to-real gap and improve downstream model performance. The modular, extensible design of the pipeline facilitates these extensions and adapts readily to related scientific plot understanding tasks.