



# Component segmentation of engineering drawings using Graph Convolutional Networks

Wentai Zhang, Joe Joseph, Yue Yin, Liuyue Xie, Tomotake Furuhashi, Soji Yamakawa, Kenji Shimada, Levent Burak Kara\*

Department of Mechanical Engineering, Carnegie Mellon University, Pittsburgh, PA, 15213, USA

## ARTICLE INFO

### Keywords:

Engineering drawings  
Graph neural networks  
Component segmentation  
Deep learning  
Computer vision

## ABSTRACT

We present a data-driven framework to automate the vectorization and machine interpretation of 2D engineering part drawings. In industrial settings, most manufacturing engineers still rely on manual reads to identify the topological and manufacturing requirements from drawings submitted by designers. The interpretation process is laborious and time-consuming, which severely inhibits the efficiency of part quotation and manufacturing tasks. While recent advances in image-based computer vision methods have demonstrated great potential in interpreting natural images through semantic segmentation approaches, the application of such methods in parsing engineering technical drawings into semantically accurate components remains a significant challenge. The severe pixel sparsity in engineering drawings also restricts the effective featurization of image-based data-driven methods. To overcome these challenges, we propose a deep learning based framework that predicts the semantic type of each vectorized component. Taking a raster image as input, we vectorize all components through thinning, stroke tracing, and cubic bezier fitting. Then a graph of such components is generated based on the connectivity between the components. Finally, a graph convolutional neural network is trained on this graph data to identify the semantic type of each component. We test our framework in the context of semantic segmentation of text, dimension and, contour components in engineering drawings. Results show that our method yields the best performance compared to recent image-based, and graph-based segmentation methods.

## 1. Introduction

Engineering technical drawings of mechanical parts serve as a universal medium for information exchange between designers and manufacturers. Such drawings encode the topological information, dimensions, and manufacturing requirements of a product in a unified and standard form, which can then be utilized in various engineering applications including content-based part indexing (Fonseca et al., 2005; Kasimov et al., 2015), cost estimation (Sajadfar and Ma, 2015), and process planning (Kulkarni et al., 2000). Although the underlying designs are commonly created in a vector format through digital design tools, a raster drawing is more frequently used by manufacturers due to the ease of information exchange and quality assurance. According to a survey of Japan's manufacturing industry (Mitsubishi UFJ Research & Consulting Co., 2019), 84% of the customers use 2D raster-based drawings such as PDF, paper, or fax format when placing an order for manufacturing, which results in a major impediment in the automation of the aforementioned applications due to the need for human involvement in interpreting these drawings.

For a modern online platform of part manufacturing, clients often upload their designs in raster image format for better quality assurance and IP protection since the information in image drawings is noneditable. Unlike a vector format, which enables trivial digital access to all stored information through a script file, raster drawings usually require manual inspection by technicians to extract the information required for quotation and manufacturing. The inspection process includes the identification of the part shape, dimensions, and manufacturing requirements.

Here, we focus on the problem of semantic segmentation of the components in raster drawings. Common mechanical engineering components consist of straight lines, arcs, and circles. Our goal is to develop an automated data-driven framework that learns to distinguish between contour shapes, dimension sets, and text at the component level (see Fig. 1). Our approach improves the efficiency of the technical drawing interpretation, relieving the burden of human operators by reducing the repetitive and tedious task of labeling drawings.

While recent vision-based methods have shown to be effective in image interpretation tasks such as object detection (Redmon et al.,

\* Corresponding author.

E-mail address: [lkara@cmu.edu](mailto:lkara@cmu.edu) (L.B. Kara).

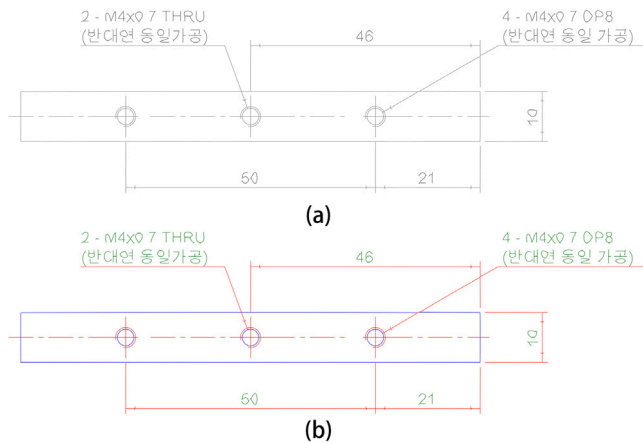


Fig. 1. (a) Typical mechanical engineering drawing. (b) Semantically labeled results. Blue: Contours, Green: Texts, Red: Dimension sets. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

2016; He et al., 2017; Wang et al., 2020), semantic segmentation (Long et al., 2015; Ronneberger et al., 2015; Chen et al., 2018), and visual question answering (Li et al., 2019; Jiang et al., 2020; Wang et al., 2021), it is challenging to apply these methods to the segmentation of components in engineering drawings since the type of a component is also dependent on the contextual information. The same component can have different semantic meanings hence categories in different contexts. Mechanical engineering drawings typically involve lines or curves only sparsely filling the image frame, and without color or textual features, making traditional image-based segmentation approaches ineffective at engineering drawing interpretation.

To address this issue, we present a new approach to engineering drawing segmentation that maps the components in the drawing to a graph structure embedded with contextual information (Fig. 2). The components are obtained through drawing preprocessing and a vectorization process, which then forms the graph nodes. The relationships between the components are encoded in the edges of the graph. A set of features are computed as nodal attributes to embed the shape, size, curvature, and neighboring information. With this, the task of component segmentation of a raster drawing image is formatted as a node-labeling problem.

Graph Convolutional Networks (GCNs) have shown great promise in node classification in various graph-structured data including academic networks (Bhattacharya and Getoor, 2007), social networks (Fonseca et al., 2005), citation networks (Sen et al., 2008) and medical data (Fakhraei et al., 2016; Namata et al., 2012). Like CNNs, GCNs aggregate information from a node and its neighboring nodes using a trainable unstructured feature map, which makes it applicable to graph data of any shape. In this work, we build our data-driven model based on GraphSAGE (Hamilton et al., 2017), a recent GCN model that aggregates the nodal information from the neighborhood structure, to predict the component type of each vectorized entity in raster drawings. The effectiveness of our graph representation is validated by the comparison between our proposed GCN methods and three vision-based or graph-based model for image segmentation. The depth of our proposed model is optimized through a parametric study of the number of convolutional layers in the model. Results indicate superior performance in both 2-class classification and 3-class classification tasks compared to our baseline models.

Our main contributions include:

- A vectorization method for raster drawings by skeletonizing, tracing, splitting, and cubic bezier curve fitting.
- A graph representation for the extracted components in engineering drawings embedded with domain-specific nodal attributes and contextual information.

- A data-driven framework that takes the vectorized component graphs as input and identifies the nodal component type for semantic segmentation.

## 2. Related works

In this section, we review some prior works for the analysis of engineering drawings. In addition, as inspiration to our proposed framework, we also introduce recent advances in graph-based methods for image analysis and general data-driven methods for graph node classification.

### 2.1. Content-based methods for engineering drawing

In the analysis of engineering drawings, content-based methods are broadly used for drawing match or retrieval. The key is to comprehend the basic elements in the drawings and define a measure of similarity. To detect the basic elements, Hough line transform (Mednionogov et al., 2000), pixel blocks (Jiao et al., 2009), and patch groups (Liu et al., 2010) are utilized as representations for extensive matching or retrieval tasks. Other works also focus on content-based detection for certain components like shape contour (Kuchuganov et al., 2020), symbols (Hu et al., 2021; Elyan et al., 2020) and information tables (Sulaiman et al., 2012) in the drawing. When an exemplar drawing is given, the matching process seeks the closest drawing in an existing pool under a similarity measure. Prior works propose distances between a graph of decomposed topology (Sousa and Fonseca, 2010), distances between feature vectors (Mednionogov et al., 2000), weighted cosine similarity (Feng et al., 2009), and the Bhattacharya correlation between histograms (Huet et al., 2001) to identify the most relevant drawing.

All works above are applied to drawings with only contour shapes, which restricts the extensive usage in our problem where the dimensions and texts are also included. But the idea of constructing a graph of basic elements enables a more efficient representation for raster drawings compared to the original image format.

### 2.2. Data-driven methods for graph node classification

When the graph of components is introduced, the problem of component segmentation is converted to the classification task of the graph nodes. Node classification is a critical problem in many supervised or semi-supervised learning scenarios (Zhu, 2005). Various methods have been proposed for node classification including iterative classification (Sen et al., 2008), label propagation (Xiaojin and Zoubin, 2002), and SVM on nodal embeddings (Grover and Leskovec, 2016). However, recent work has shown the promise of a better classification if the nodal embedding is jointly learned with the training of the data-driven classifier (Yang et al., 2016), which drives the development of end-to-end graph neural networks.

Since the first deep learning based framework, Deepwalk (Perozzi et al., 2014), was introduced to solve the nodal classification problem, graph neural networks have demonstrated their superior capability in efficient feature extraction on unstructured graph data. Recently, Graph Convolutional Networks (GCNs) (Welling and Kipf, 2016) have yielded better performance due to the unique nodal information aggregation mechanism. Other variants of GCN extensively introduce various ways for aggregation including neighborhood aggregation (Hamilton et al., 2017) and attention mechanism (Veličković et al., 2017), which inspires us to build our proposed data-driven model for component segmentation.

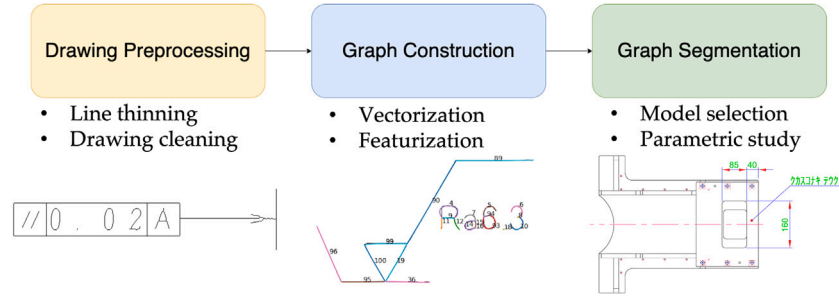


Fig. 2. Our proposed workflow.

### 2.3. Graph-based methods for image analysis

For a more efficient and structured feature representation, graph models have been frequently introduced to multiple major tasks in image analysis. Monti et al. (2017) propose the first GNN-based method in image classification. More recent works focus on achieving image segmentation using the graph of pixels (Shi and Malik, 2000) or graph of superpixels (Stutz et al., 2018). These oversegmented, simplified, images can be applied to many classic tasks in computer vision, including depth estimation, segmentation, and object localization (Achanta et al., 2012).

The biggest advantage of introducing the graph representation in image analysis is that the problem can be decomposed into a higher level image processing with a more coarse unstructured data format. Unlike the image space where each pixel is treated as an independent element, the connection between the graph nodes enforces the possible dependencies due to the continuity of images. Paliwal et al. (2021) proposed a symbol detection method based on Graph Convolutional Networks for piping diagrams. The one-shot learning model yielded comparable performance to previous fully supervised model. In a more recent work, Rica et al. (2023) introduced a zero-error digitization approach for piping diagrams by creating a graph of symbols and components. The proposed tools were able to indicate incorrectly identified components to aid manual validation as well as search groups of components based on a given query. In both works, graph-based feature extraction methods have demonstrated their superior performance in parsing contextual information in structured diagrams. Despite the similarity of these works to our work, their primary focus is on symbol recognition, while the focus of our proposed method is on component node classification. In mechanical engineering drawings, the type of a component is heavily dependent on the contextual information and the neighboring components. As an initial attempt to apply graph-based methods in the analysis of mechanical parts, Xie et al. (2022) developed a data-driven framework that can identify the manufacturing process of a part using a graph of detected straight-line segments in the engineering drawing. But the graph is not suitable for a complete semantic interpretation of all the components. First, only straight-line segments are used to vectorize the drawing, which leads to inaccurate vectors for circles, arcs and curved strokes in text. Second, The featurization of the vectors only contains location information (X, Y coordinates). There are no indicators to describe the topological feature (size, angle, curvature) of each obtained vector. Therefore, this work is shown to be effective in graph classification tasks (manufacturing method classification), but it is challenging to be directly extended to graph segmentation tasks (component interpretation).

As such, we propose a novel method to vectorize the drawing as a lower-level component representation with lines and curves, and construct a graph of such vectors as the basic element with topological featurization for drawing analysis, embedding the contextual information in the edges of our component graph. The proposed graph representation is utilized to achieve a component segmentation of engineering drawings.

### 3. Technical approach

In this work, we present a pipeline to preprocess engineering drawings, construct a labeled drawing dataset and train a graph-based neural network model for component segmentation. This pipeline starts with a method that converts the raster drawing images into vectorized curves. Subsequently, a self-defined component graph is constructed based on the connections and distances among the obtained components. For each node (component) in the graph, we also present a novel featurization method to generate a series of feature parameters based on sampled points. Finally, this feature-embedded graph is utilized as input to a graph convolutional neural network that is able to predict the component type of each node (vector) in the drawing. Our inference pipeline is summarized in Algorithm 1

---

#### Algorithm 1: Engineering Drawing Graph Network (EDGNet)

---

**Input :** A raster engineering drawing  $D_r$

**Output:** A graph of vectorized components with semantic labeling  $G(N, E), Y$

---

- 1  $D_s \leftarrow \text{skeletonize}(D_r) + \text{smoothing};$
  - 2 Split  $D_s$  to strokes  $\{S_i\}_{n_N}$  from junction points;
  - 3 Fit cubic bezier curves  $\{B_i\}_{n_N}$  to each stroke  $S_i$ ;
  - 4 Init  $i = 1$ ;
  - 5 **while**  $i \leq n_N$  **do**
  - 6     Sample  $n$  equally spaced points  $P_i$  on  $B_i$ ;
  - 7     Calculate nodal features  $N \in \mathbb{R}^{n_N \times (4n-1)}$ ;
  - 8      $i = i + 1$ ;
  - 9 **end**
  - 10 Assemble a component graph  $G(N, E)$ ,  $E$  are the edges (connections) between the nodes;
  - 11 Get predicted semantic labels from a graph convolution network  $Y \in \mathbb{R}_N^n \leftarrow F_\theta(G(N, E))$
- 

#### 3.1. Drawing vectorization

As mentioned in Section 1, vision-based methods are not effective for feature extraction from raster engineering drawings since the information exists sparsely as discrete black pixels. Engineering Drawings which include sheet metal parts, lathing parts, and general machining parts in their raster format do not contain any component type information and are greyscale. Convolutional neural networks have severe limitations in embedding distant contextual information on very large images due to their ordered grid structure (Simonyan and Zisserman, 2014) and require endless increasing in depth. However, engineering drawings in their original vectorized form are capable of encoding long-scale connectivity through an arbitrary graph structure. The basic elements of this graph structure consist of lines, curves, and texts of nodes and connectivity between them. As such, we propose a novel method to vectorize the raster drawings before analysis. The method consists of three broad steps: skeletonization, trajectory tracing, and curve fitting.

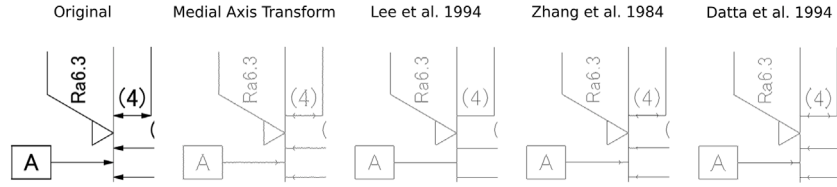


Fig. 3. Comparison of four thinning methods in skeletonizing engineering drawings.

A CAD drawing is a collection of parametrically stored entities consisting of headers, blocks, tables, entities, and objects that can be easily rendered into a raster drawing through lines, curves, and text. However, it is challenging to identify these parameters from raster drawings due to the width of each entity and binary colors in the rendering process. It is also challenging to vectorize entities that are overlapping or lines that are intersecting, as the vectorization input only receives the rendering of the engineering drawing and not the originally created file from the CAD or drafting tool. As such, it is necessary to morphologically thin the drawing to identify the central trace of these strokes in raster drawings and extract the skeleton of each line for better fitting. In our implementation, three line thinning algorithms (Zhang and Suen, 1984; Lee et al., 1994; Datta and Parui, 1994) were tested for extracting the skeleton out of the original drawing (Fig. 3). From the tests, it was concluded that the Medial Axis Transform method tends to generate more small branches where the line width is relatively large. The method from Lee's cannot retain a complete structure for small components like the arrowheads. Zhang's method does not produce clean junction points. Datta et al. method was selected to have the most desirable thinned morphology for further processing and generating the parametric curves.

1. Thinning or skeletonization of the image to form traces with a single pixel width.
2. Smoothing the pixels to obtain a set of traces.
3. Splitting the traces to account for corners within the trace.
4. Removing small traces, and merging junction points
5. Fitting cubic bezier curves

Through the skeletonization process, strokes of lines and curves in the raster drawing are thinned to single pixel-wide trajectories. This allows us to trace the trajectories through neighboring pixels and convert the entire drawing into a unified set of parametric curves for better segmentation. In Fig. 4, we summarize three types of points existing in the obtained skeleton. The point type is defined based on the number of black pixels it connects to, which is efficiently calculated by a filter scanning over all the black pixels (Fig. 4(a)).

We define a trace as an ordered list of connected pixels starting from either a junction point or endpoint and ending at either. To obtain this set of traces, we start at a randomly selected black pixel (Fig. 4(b)) and iteratively visit its neighboring points until a termination point (a junction or endpoint) has been reached. We then reverse the trace to reach the other termination point. This trajectory is recorded, forming a completed trace used for parametric curve fitting. We continue this process until all black pixels have been visited and an entire set of traces have been created.

Now, with the set of pixel traces from the image, we split the traces if an edge is detected within the trace. To break the trace, we evaluated the cosine angle of the vectors formed between the point on the trace  $p_i$  to the two terminal points;  $p_s$  and  $p_e$ , given by;

$$\theta_i = \cos^{-1}(\overrightarrow{p_s - p_0} \cdot \overrightarrow{p_e - p_0} / |\overrightarrow{p_s - p_0}| \cdot |\overrightarrow{p_e - p_0}|)$$

Taking the second derivative of this angle provides a clear indication of a corner within the trace by forming a spike in the second derivative. We find these spikes in the  $\bar{\theta}$  vector by finding the local maxima by a simple comparison of neighboring values. The trace is split at this spike point  $p_i$ . All the traces are split based on these criteria.

Subsequently, we check for small traces and eliminate them. Traces with  $m_j < 4$  pixels are eliminated to ensure that we do not have an underdefined problem while fitting the cubic bezier curves in the next step. We also merge the terminal points of small traces, so the neighboring traces are now connected. We also merge junction points to ensure the edge-connectivity of the graph. Finally, we fit cubic bezier curves on traces. We use the least square method for fitting the cubic Bezier curves as shown in the equation:

$$\bar{p}^j = [p_0^j, p_1^j, \dots, p_{m_j-1}^j]^T$$

$$B_{i,n_{ord}}^j(t) = \binom{n_{ord}}{i} (t)^i (1-t)^{n_{ord}-i}, i = 0, 1, \dots, n_{ord}, t \in [0, 1]$$

$$B_{m_j \times (n_{ord}+1)}^j x_{(n_{ord}+1) \times 2}^j = \bar{p}_{m_j \times 2}^j$$

where,  $\bar{p}^j$  is the ordered list of  $m_j$  points in the  $j$ th trace,  $B_{i,n_{ord}}^j(t)$  is the Bernstein polynomial of order  $n_{ord} = 3$ ,  $B^j$  is the Bernstein matrix of the cubic bezier curve and  $x^j$  is the control points for the  $j$ th trace. We set the terminal points as the first and last points of the bezier curve control points, i.e.  $x_0^j = p_s^j$  and  $x_3^j = p_e^j$  (see Figs. 5 and 6).

These steps provide the necessary set of parameterized vectors needed for the construction of a unified graph used in our approach for segmenting the engineering drawing into its various components.

### 3.2. Graph construction

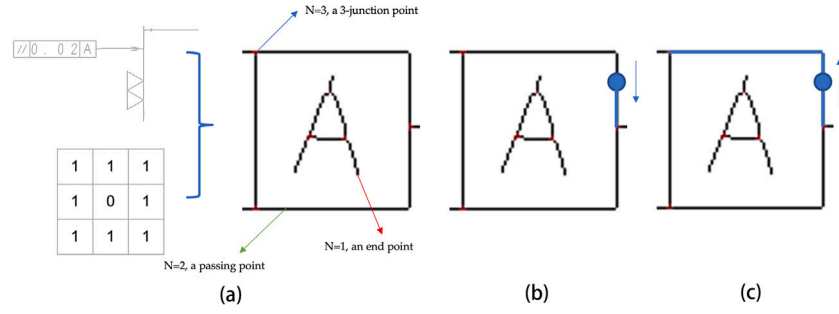
At the heart of our proposed method, a graph representation for the components is the key to embedding topological features and contextual information among the distant but connected components. In our graph structure, the node vectors are obtained from our vectorization. The graph edges are generated from the connection between these vectors. To form a unified representation, we sample  $n$  evenly spaced points along each vector for featurization. Table 1 lists the features computed based on these sampled points. The proposed features encode the shape, length, angle, and curvature information from each vector. Note that the entire drawing is normalized to fit in a unit square before being fed into the featurization process to ensure that all features listed are independent of the input drawing size. In summary, our graph model is defined as:

$$G(N, E), N \in \mathbb{R}^{n_N \times (5n-1)}, E \in \mathbb{Z}^{n_E \times 2}$$

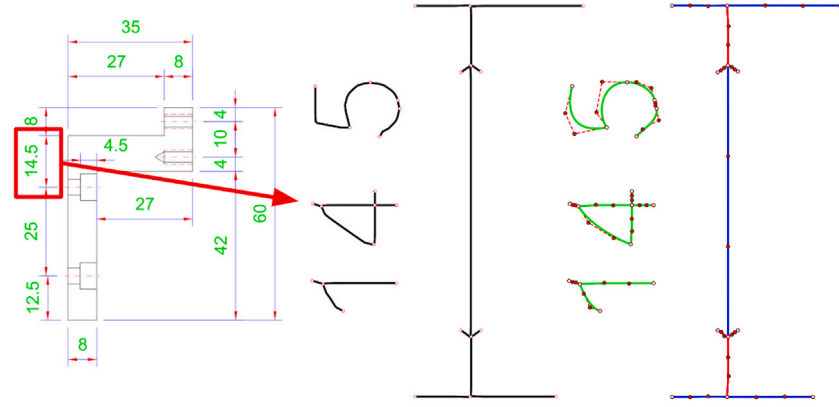
where,  $N$  are the nodal features of a graph with  $n_N$  nodes,  $E$  are the edge indices of a graph with  $n_E$  edges. In this way, the target task of this work is converted to predicting the nodal labels. The label, as well as the component type, is defined as  $Y \in \mathbb{Z}_N^n$ . It is a categorical parameter list indicating the type (contour lines, dimension lines, or texts) of each vector.

To establish the ground truth labels  $Y_{gt}$  when creating a dataset for training, we modify a batch renderer for DXF drawings. The renderer EZDXF (<https://ezdxf.mozman.at/>) is modified to be able to parse the component type information stored in the vector DXF drawing and paints each component type with a unique color (Fig. 7). For each vector in the graph, we sample points and check the corresponding color of each point location in the ground truth image. Finally, the ground truth label for each vector is determined by majority voting

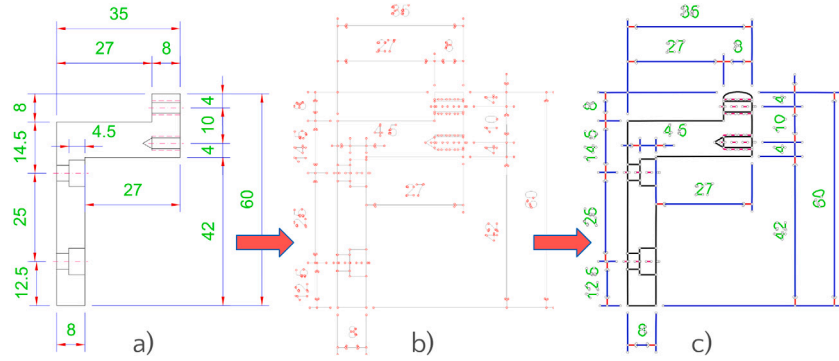




**Fig. 4.** (a) The process of identifying the junction points, passing points and end points of a skeletonized drawing. (b) The trajectory tracing process from a passing point to a junction or an end point. (c) After (b), another tracing is done to search for the other end of the trajectory.



**Fig. 5.** Raster drawing with the highlighted region showing the thinned pixels with junction points and subsequently the cubic bezier curves fitted on the traces with their respective control points.



**Fig. 6.** (a) Raster drawing (b) Thinned, smoothed, and traced trajectories with junction points marked in red circles (c) Cubic bezier curves fitted on pixel traces.

**Table 1**

Our proposed nodal features.  $n$  is the number of sampled points on each vector.

Feature type	Parameter	Dimension
Shape	$n$ points sampled along the trace, XY coordinate	$2n$
Length	Length between each pair of consecutive points	$n - 1$
	Total length (of the curve)	1
	First-to-last/total	1
Angle	Cos angle between each pair of consecutive short line segments	$n - 2$
Curvature	Curvature at each sampled point	$n$

of the sampled points. In this work, we conduct two testing conditions in terms of the segmentation labels: (1) text/non-text: The

model is designed to distinguish the green components vs the other ones. (2) text/contour/dimension: The model is designed to distinguish the green, the black, and the other components. The labels  $Y_{gt}$  are converted to one-hot encoding accordingly.

### 3.3. Graph segmentation

Based on the graph representation explained above, our task can be represented as:

$$Y \leftarrow F_{\theta}(G(N, E))$$

where  $G(N, E)$  is our proposed graph structure,  $F_{\theta}$  is a data-driven model with trainable parameters  $\theta$ ,  $Y$  is the predicted component type label from the model. During the training of  $F_{\theta}(\cdot)$ ,  $Y_{gt}$  is utilized for loss calculation. Then, we introduce a model  $F_{\theta}(\cdot)$ , a loss function  $L$ , and an optimizer to launch the training process.

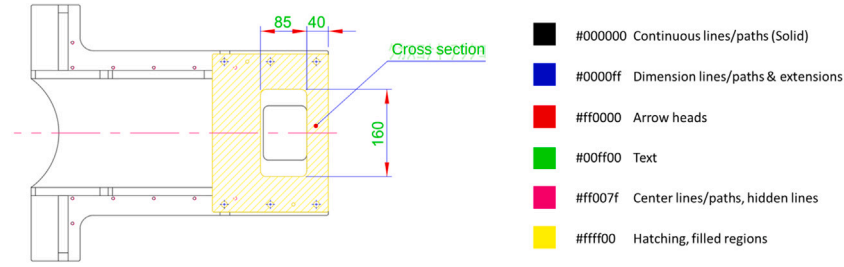


Fig. 7. A sample rendering of a DXF drawing used for ground truth label retrieval. Components of each type are painted with a unique color. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

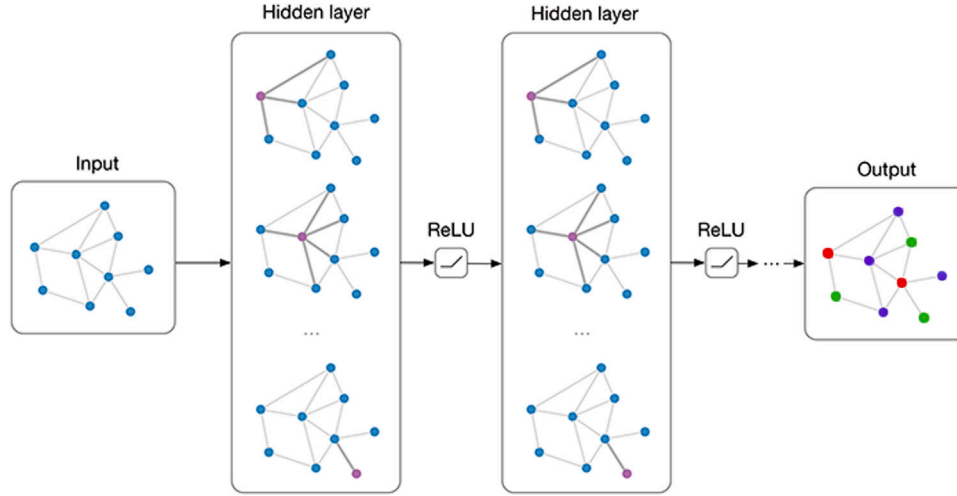


Fig. 8. A general structure of Graph Convolutional Network (GCN) in nodal classification tasks.

Graph Convolutional Networks (GCNs) (Fig. 8(a)) is a type of network that works on arbitrarily structured graphs. They have been shown to be effective in various nodal classification tasks including citation and social networks (Welling and Kipf, 2016; Hamilton et al., 2017). GraphSAGE model shows superior performance in predicting the nodal class label over the other baseline graph models (Hamilton et al., 2017). Therefore, we build our model based on GraphSAGE. A GCN (Welling and Kipf, 2016) model and an MLP model are also implemented to serve as baseline models. As a parametric study, we vary the number of convolutional layers in the experiments for the optimized network architecture.

Based on the steps explained above, a graph dataset is constructed based on 430 real engineering part drawings, including sheet metal parts, lathing parts, and general machining parts, from a large e-commerce system for custom mechanical parts. The drawings are in DXF format originally and converted to black-and-white raster images as training data. The number of vectorized components in each drawing ranges from 500 to 3000. The dataset is split into 80/20 for training and validation. During training, the loss function is chosen as the cross-entropy loss between the predicted nodal class and the ground truth label. Adam (Kingma and Ba, 2014) optimizer is utilized with a learning rate  $1e-3$ , weight decay  $5e-4$ . All the models are trained with a maximum epoch of 10,000 and batch size of 16. The best model regarding the validation accuracy is saved for inference. In our experiments, each model takes about 20 h to train on a GeForce RTX 2080 TI Graphics Card.

#### 4. Results

This section demonstrates a series of experiments to validate the effectiveness of our graph representation. Additionally, Experiment II focuses on optimizing the model architecture with respect to the

validation accuracy. Finally, this model setup is extended to a 3-class segmentation problem.

##### 4.1. Model selection

Using our graph representation described in Section 3.2, a classifier is trained to predict the component type of each vector in the drawing. Here, we introduce three data-driven classifiers, including our proposed GraphSAGE model (GS), a vanilla GCN model (GCN) as a graph method baseline, and a Multi-layer Perceptron model (MLP) as a non-graph method baseline in our designed Experiment I. To maintain a fair comparison, all three models are designed to have similar architecture and depth. The details are:

- **GCN model:** 3 vanilla graph convolutional layers + 2 linear layers, number of nodes in each hidden layer: [32, 64, 128, 32].  $4 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.
- **GS model:** 3 GraphSAGE convolutional layers + 2 linear layers, number of nodes in each hidden layer: [32, 64, 128, 32].  $4 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.
- **MLP model:** 5 linear layers, number of nodes in each hidden layer: [32, 64, 128, 32].  $4 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.

In Experiment I, the three models described above are trained with identical training conditions to distinguish the text vs non-text components in our constructed part drawing dataset. We use  $n = 4$  for the dataset creation in this experiment based on a parametric study detailed in Appendix A. Fig. 9 illustrates the validation accuracy in the training process. It can be concluded that graph-based models (GS and GCN) yield better ( $> 5\%$ ) results than the non-graph-based model (MLP), which speaks for the necessity of the contextual information embedded in our designed graph structure. Conversely, GS and GCN

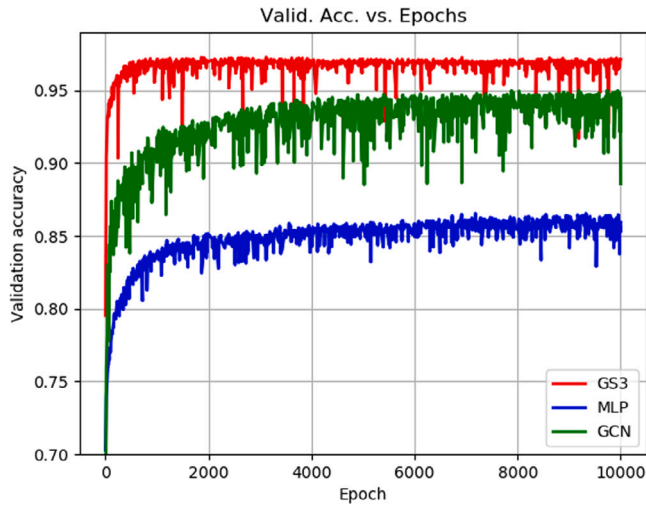


Fig. 9. The training curves for GS, GCN and MLP model in Experiment I.

model share the same pattern in the validation curve in the early phase of the training ( $< 500$  epochs). Then, the GS model rapidly converges at around 1000 epochs to 95%, while the GCN model gradually converges to a lower value at around 3000 epochs. The results confirm the superior capability of our implemented GS model. Next, we want to boost its performance by optimizing the model depth.

#### 4.2. Model depth study

Through Experiment I, the GS model with 3 convolutional layers yields the best performance over the other two models. Therefore, we further design Experiment II to study the effect of model depth on the final classification accuracy. Three GS models in this experiment are implemented with 3, 4, and 5 graph convolutional layers. For consistency, all three networks are assembled so that the graph layers expand the dimensions. Then the linear layers squeeze the dimensions to the final desired number of classes. The details are:

- **GS3** model: 3 GraphSAGE convolutional layers + 2 linear layers, number of nodes in each hidden layer:  $[32, 64, 128, 32]$ .  $4 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.
- **GS4** model: 4 GraphSAGE convolutional layers + 3 linear layers, number of nodes in each hidden layer:  $[32, 64, 128, 256, 128, 32]$ .  $6 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.
- **GS5** model: 5 GraphSAGE convolutional layers + 4 linear layers, number of nodes in each hidden layer:  $[32, 64, 128, 256, 512, 256, 128, 32]$ .  $8 \times \text{ReLU} + 1 \times \text{Softmax}$  as nonlinear activation.

We utilize the same setup as Experiment I for training all three models above. Fig. 10 demonstrates the validation curves for these GS models during the training session. We conclude that the performance of the GS model tends to increase for deeper models. However, the improvement gradually levels out when 5 convolutional layers are used. Additionally, Table 2 summarizes the confusion matrix for the best model, GS5. Sample prediction results are also demonstrated in Fig. 11. Results show that our GS5 model achieves nearly perfect results in all three test drawings. More statistical comparison results with other baseline models are detailed later in Section 4.4. Then, we extend the test condition to multi-class segmentation.

#### 4.3. Multi-class segmentation

In previous experiments, the classifier is trained to distinguish between the text and non-text components. This task is relatively

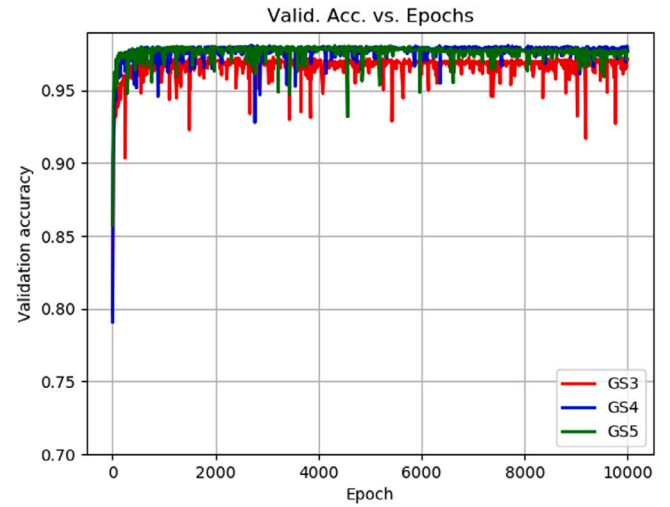


Fig. 10. The training curves for GS models of different depths in experiment II.

Table 2

Confusion matrix of our proposed GS5 model in the 2-class segmentation task. From these statistics, the evaluation results can be calculated: **Precision** = 98.38%, **Recall** = 98.57%, **Accuracy** = 98.48%.

GT	Prediction	
	Text	Non-text
Text	21 103	304
Non-text	345	13 559

Table 3

Confusion matrix of our proposed GS5 model in the 3-class segmentation task. Based on the statistics, the performance of the model can be calculated: **Precision**: 83.03%, 95.04%, 90.70%, **Recall**: 83.46%, 96.07%, 89.60%, **Accuracy**: 90.82%.

GT	Prediction		
	Contour	Text	Dimension
Contour	4238	130	710
Text	105	9229	273
Dimension	761	352	9589

simple since text components are usually unique in terms of size and curvature compared to all other straight lines and curves. For practical use, the human inspector needs to comprehend the overall shape of the part through all contour lines, and then gather all manufacturing requirements through dimensions and texts. As such, we construct a dataset with three component types for prediction: Contour, Text, and Dimension, which correspond to black, green, and all other colored lines in Fig. 7. The dataset is utilized in Experiment III for training a GS5 model in a 3-class segmentation task for the part drawings.

The resulting confusion matrix of our proposed model (GS5) is shown in Table 3. It can be concluded that the model retains high accuracy on the text components, while the separation between the contour lines and dimension sets is more challenging to learn. Fig. 12 illustrates some failure cases in the prediction results. Two major types of failure cases are: (1) Isolated small components in the text. For example, misclassification happens on the dashed line, diameter symbol, and through hole symbols in the sample results. A potential cause is that these components are not connected with any other components in the drawing, which makes the prediction equivalent to judging the component type only by its topological features without contextual information. The issue can be resolved if we also generate edges for the nearest neighbors of each component. (2) The region where different types meet, like the contour line with its correlated extension line. The issue is likely a result of the lack of information on the connections between components. There is no indication of the difference between

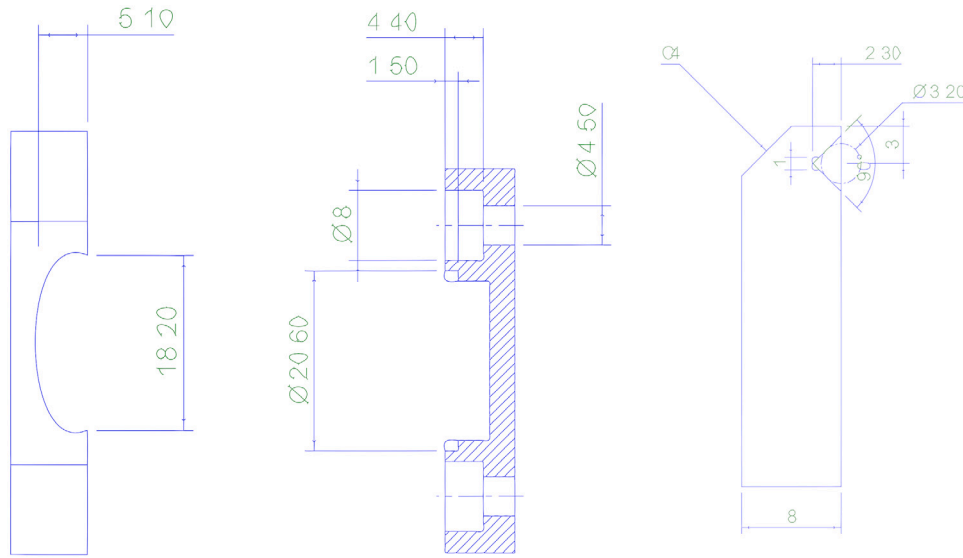


Fig. 11. Sample prediction results from GS5 models in experiment II. The colors green and blue indicate the predicted text components and non-text components respectively. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

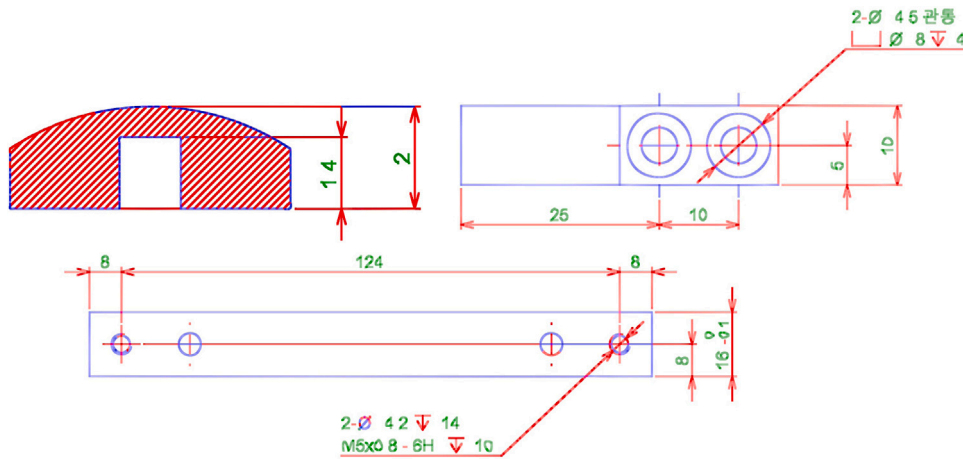


Fig. 12. Sample prediction results from GS5 in Experiment III. Blue, green and red indicates the predicted contour components, text components, and dimension components respectively. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

a line-line connection and a line-curve connection. In our graph design, a series of edge features should also be added to provide such insights into the graph network model.

#### 4.4. Baseline comparison

To better understand the challenge, we construct three baseline models to achieve the component segmentation tasks with the same dataset, including two image-based deep learning models based on PSPNet (Zhao et al., 2017) and DeepLabV3 (Chen et al., 2017), and one graph-based method based on Sketchgmn (Yang et al., 2021). For image-based models, a black-and-white drawing image is fed to the model to predict a map of indices to indicate the semantic label of each pixel. To transfer the pixel level prediction to the component level, we map the predicted labels to vectorized results using majority voting from all pixels each vector is passing.

Using the same set of data for training and validation, the resulting validation accuracy for all three baseline models and our model are summarized in Table 4. Here, we compare the models in three segmentation tasks: (1) Text vs. Non-text (Contour+Dimension). (2) Contour vs. Non-contour (Text+Dimension). (3) Text vs. Contour vs. Dimension.

Table 4

Comparison results for accuracy of the predicted component labels on the validation set (%) in three segmentation tasks.

Validation accuracy % Task	PSPNet	DeepLabV3	Sketchgmn	EDGNet (Ours)
Text/Non-text	96.62	96.87	94.76	<b>98.48</b>
Contour/Non-contour	80.54	82.57	88.10	<b>94.57</b>
Text/Contour/Dimension	79.54	81.64	84.37	<b>90.82</b>

It can be concluded that our model yields the best performance in all three tasks. Evident improvement can be seen in the separation between contour and dimension (tasks 2 and 3) when using graph-based models, which reinforces the idea that it is more challenging for image-based approaches to parsing sparse, man-made images. From the visual comparison shown in Fig. 13, it can be seen that our model is the only one that successfully identifies both the hole and thread line in the first example. Additionally, the other models usually have difficulty when there are multiple concentric circles with center lines. These misclassified components can easily mislead the model when extracting the overall shape of the part. More comparison results are demonstrated in Appendix B.



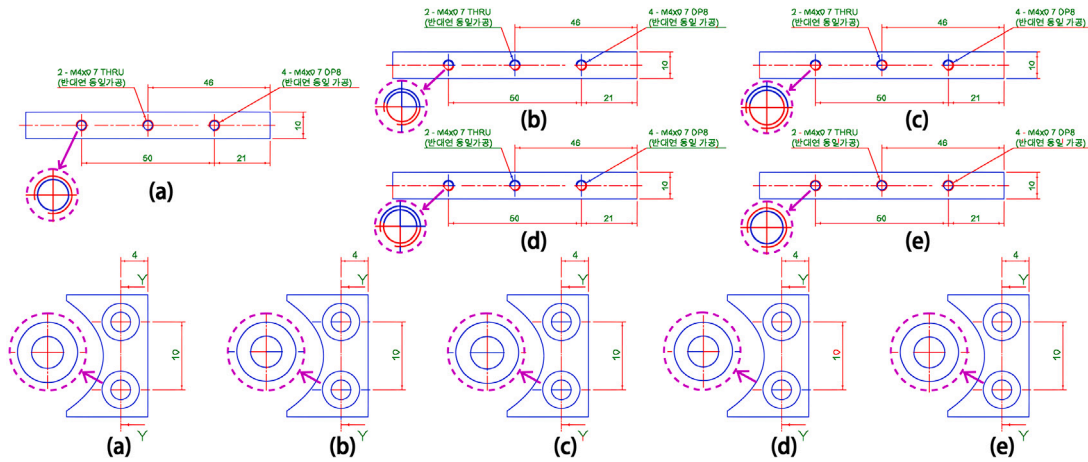


Fig. 13. Sample prediction results from three baseline models versus ours. (a) The ground truth. (b) PSPNet results. (c) DeepLabV3 results. (d) Sketchgnn results (e) Ours.

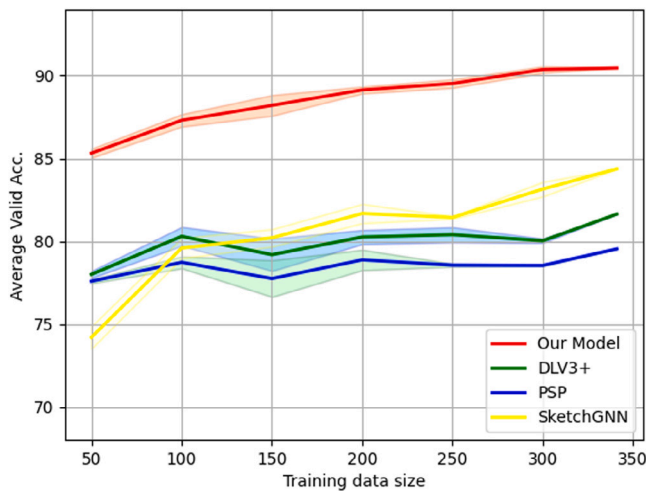


Fig. 14. Validation accuracy of the models trained on datasets of various sizes.

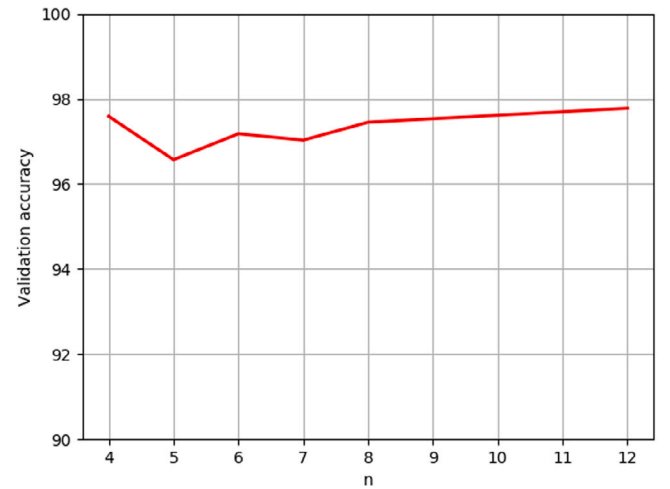


Fig. A.15. Validation accuracy of GS model trained on datasets with the number of sampled points  $n$  increases.

To further compare the model stability when applied to datasets of various sizes, we conduct another set of experiments using a subset of our current train data. Each time,  $k \in [50, 340]$  drawings are randomly sampled from our training set for training all four models in the task of 3-class segmentation. Then we repeat the process five times with different random seeds to eliminate the influence of a biased subset taken by accident. The resulting validation accuracy curves for all four models are demonstrated in Fig. 14. It can be concluded that overall our model outperforms the other baseline models consistently by over 5%. Another interesting finding is that graph-based methods tend to have less variance when different subsets are used compared to image-based methods, which illustrates better stability of feature extraction when encountering different data.

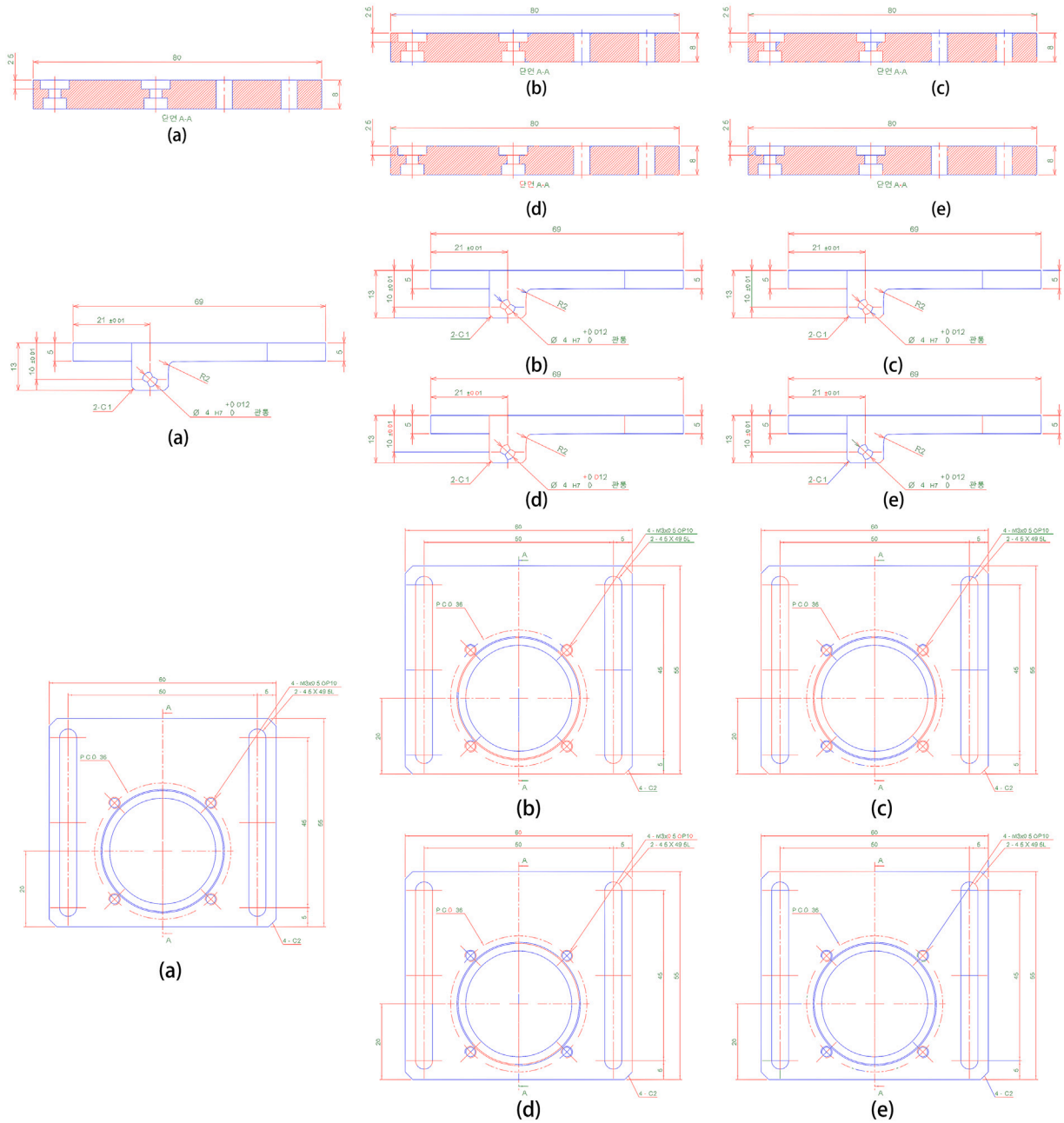
## 5. Discussions and future work

With our current approach, we are able to automate the vectorization and labeling process of raster engineering drawings in terms of the separation among contour lines, dimension sets, and texts. For broader practical use in the industry, symbols such as manufacturing requirement symbols and geometric tolerance symbols are also critical in the quotation process. Considering they are usually consistent in shape and style, an immediate next step of our current work is to develop a detection algorithm to identify such symbols in the obtained vectorized results from our preprocessing step. Simple heuristic-based

methods can be applied to searching for surface roughness symbols or through hole symbols since they retain highly consistent shapes (equilateral triangles) across all types of drawings. Data-driven models should be utilized to locate more complex symbols like geometric tolerance symbols because they are usually a composite of texts, symbols, and indication boxes.

From our results on three-class segmentation, it can be concluded that contextual information about the connection between two components is needed for better accuracy. For example, edge features to indicate the angle, the shift and the curvature change between the two connected components can be added to our current graph representation. This information can help the network to understand if two components are truly connected based on semantic meaning or just independent of each other with an intersection. With this new graph representation, GraphSAGE model should also be replaced with more advanced graph networks that also take edge attributes as input for analysis, such as Graph Attention Networks (Veličković et al., 2017), Graph Transformers (Dwivedi and Bresson, 2020), and GINE Convolution (Hu et al., 2019).

The ultimate goal of our work is to aid a human operator when inspecting drawings for topology and manufacturing information necessary for the later quotation process. The effectiveness of our framework needs to be practically validated by human users. To enable easier access for the system we propose in this work, an interactive user



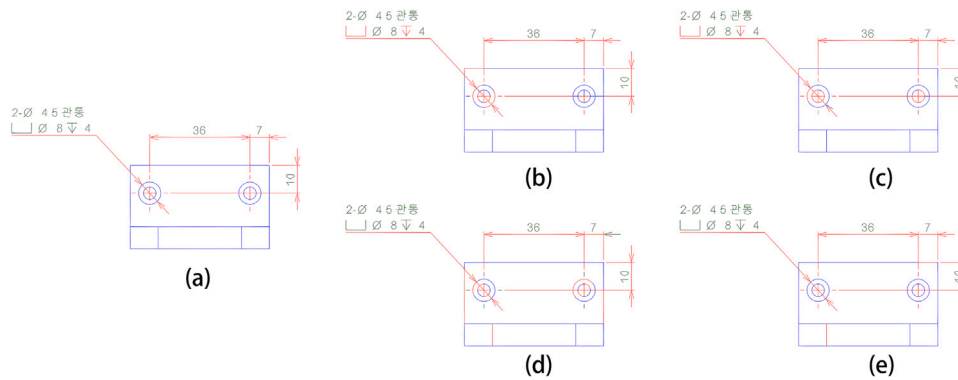
**Fig. B.16.** Sample prediction results from three baseline models versus ours. (a) The ground truth. (b) PSPNet results. (c) DeepLabV3 results. (d) Sketchgnn results (e) Ours. In the first part, our model yields the best results when identifying all the internal surfaces of the holes as contour lines (blue). In the second and third part, our model is the only one that correctly identifies the outline of holes as blue and the center lines on the holes as red. This is crucial when extracting the shape of the entire part. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

interface should be developed to allow the users to upload their own drawings, launch the vectorization, and get the automatic component type prediction results. Then the user is only responsible for marking the critical information or correcting minor errors in the prediction. Compared to the original inspection and labeling task, the work for the human operator is more efficient and intellectual.

## 6. Conclusions

In this work, we present a novel framework for raster engineering drawing analysis, including a preprocessing method for drawing vectorization, a graph representation embedded with domain knowledge,

and a data-driven model to learn and predict the component type of each vector in the drawing. The framework converts the problem from sparse image comprehension into semantic segmentation of the vectorized components from the original drawing, enhancing the efficiency of feature extraction. Results also show that our method yields superior performance in distinguishing the semantic meaning of contour/dimension lines compared to common CNN-based image segmentation methods. A similar framework can be established to other analyses of raster engineering drawings such as manufacturing method classification (Xie et al., 2022), dimension estimation, and similarity search. The proposed graph representation has the potential to be used extensively in developing a digitized tool for a part quotation.



**Fig. B.17.** Sample prediction results from three baseline models versus ours. (a) The ground truth. (b) PSPNet results. (c) DeepLabV3 results. (d) Sketchgnn results (e) Ours. In this part drawing, our model is that only one that correctly distinguishes the outline of the holes as blue and the center lines on the holes as red. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

### CRedit authorship contribution statement

**Wentai Zhang:** Conceptualization, Methodology, Writing – original draft, Software, Data curation, Visualization. **Joe Joseph:** Writing – original draft, Software, Data curation, Visualization. **Yue Yin:** Software, Data curation, Visualization. **Liuyue Xie:** Investigation, Writing – review & editing. **Tomotake Furuhashi:** Conceptualization, Writing – review & editing. **Soji Yamakawa:** Conceptualization, Writing – review & editing. **Kenji Shimada:** Project administration, Funding acquisition, Writing – review & editing. **Levent Burak Kara:** Project administration, Funding acquisition, Writing – review & editing.

### Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Wentai Zhang, Joe Joseph, Quan Chen, Liuyue Xie, Amit Regmi, Soji Yamakawa, Tomotake Furuhashi, Kenji Shimada and Levent Burak Kara reports financial support was provided by Misumi Group Inc.

### Data availability

The authors do not have permission to share data.

### Acknowledgments

The authors would like to thank MiSUMi Corporation for their provision of a contemporary engineering problem, guidance on the applicability of developed methods, and financial support. Additionally, the authors appreciate the help from Run Wang, Zhuoran Cheng, Zheren Zhu, and Chenlai Wang.

### Appendix A. Parametric study on $n$

In this study, we vary the number of sampled points on each vector  $n \in \{4, 5, 6, 7, 8, 10, 12\}$  to explore its effect on the final model performance. Like GS3 model detailed in Section 4, as  $n$  increases, we enlarge the width of each layer accordingly. The resulting validation accuracies as  $n$  increases are summarized in Fig. A.15. From the figure, it can be concluded that there are no significant changes in validation accuracy as  $n$  varies. A potential cause lies in the fact that the majority of the obtained vectors are straight lines. There is no extra useful information added to the input when more points are sampled in between. The difference when  $n = 4$  and  $n = 12$  can be ignored. But  $n = 12$  requires 3x more parameters to train, which usually leads to much more training time and less stability. As such, we choose  $n = 4$  for all of the experiments for feature extraction.

### Appendix B. More results on 3-class segmentation

As a supplement to the baseline comparison results shown in Section 4.4, more visual comparisons are demonstrated in Figs. B.16 and B.17.

### References

- Achanta, R., Shaji, A., Smith, K., Lucchi, A., Fua, P., Süsstrunk, S., 2012. SLIC superpixels compared to state-of-the-art superpixel methods. *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (11), 2274–2282.
- Bhattacharya, I., Getoor, L., 2007. Collective entity resolution in relational data. *ACM Trans. Knowl. Discovery Data (TKDD)* 1 (1), 5–es.
- Chen, L.-C., Papandreou, G., Schroff, F., Adam, H., 2017. Rethinking atrous convolution for semantic image segmentation. *arXiv preprint arXiv:1706.05587*.
- Chen, L.-C., Zhu, Y., Papandreou, G., Schroff, F., Adam, H., 2018. Encoder-decoder with atrous separable convolution for semantic image segmentation. In: *Proceedings of the European Conference on Computer Vision. ECCV*, pp. 801–818.
- Datta, A., Parui, S.K., 1994. A robust parallel thinning algorithm for binary images. *Pattern Recognit.* 27 (9), 1181–1192.
- Dwivedi, V.P., Bresson, X., 2020. A generalization of transformer networks to graphs. *arXiv preprint arXiv:2012.09699*.
- Elyan, E., Jamieson, L., Ali-Gombe, A., 2020. Deep learning for symbols detection and classification in engineering drawings. *Neural Netw.* 129, 91–102.
- Fakhraei, S., Sridhar, D., Pujara, J., Getoor, L., 2016. Adaptive neighborhood graph construction for inference in multi-relational networks. *arXiv preprint arXiv:1607.00474*.
- Feng, G., Viard-Gaudin, C., Sun, Z., 2009. On-line hand-drawn electric circuit diagram recognition using 2D dynamic programming. *Pattern Recognit.* 42 (12), 3215–3223.
- Fonseca, M.J., Ferreira, A., Jorge, J.A., 2005. Content-based retrieval of technical drawings. *Int. J. Comput. Appl. Technol.* 23 (2–4), 86–100.
- Grover, A., Leskovec, J., 2016. node2vec: Scalable feature learning for networks. In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 855–864.
- Hamilton, W., Ying, Z., Leskovec, J., 2017. Inductive representation learning on large graphs. *Adv. Neural Inf. Process. Syst.* 30.
- He, K., Gkioxari, G., Dollár, P., Girshick, R., 2017. Mask r-cnn. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 2961–2969.
- Hu, W., Liu, B., Gomes, J., Zitnik, M., Liang, P., Pande, V., Leskovec, J., 2019. Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*.
- Hu, H., Zhang, C., Liang, Y., 2021. Detection of surface roughness of mechanical drawings with deep learning. *J. Mech. Sci. Technol.* 35, 5541–5549.
- Huet, B., Kern, N.J., Guarascio, G., Meriardo, B., 2001. Relational skeletons for retrieval in patent drawings. In: *Proceedings 2001 International Conference on Image Processing (Cat. No. 01CH37205)*, Vol. 2. IEEE, pp. 737–740.
- Jiang, H., Misra, I., Rohrbach, M., Learned-Miller, E., Chen, X., 2020. In defense of grid features for visual question answering. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10267–10276.
- Jiao, L., Huang, F., Teng, Z., 2009. An engineering drawings retrieval method based on density feature and improved moment invariants. In: *Proceedings of the 2009 International Symposium on Information Processing (ISIP'09)*.
- Kasimov, D.R., Kuchuganov, A.V., Kuchuganov, V.N., 2015. Individual strategies in the tasks of graphical retrieval of technical drawings. *J. Vis. Lang. Comput.* 28, 134–146.
- Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

- Kuchuganov, V.N., Kuchuganov, A.V., Kasimov, D.R., 2020. Clustering algorithm for a set of machine parts on the basis of engineering drawings. *Program. Comput. Softw.* 46, 25–34.
- Kulkarni, P., Marsan, A., Dutta, D., 2000. A review of process planning techniques in layered manufacturing. *Rapid Prototyp. J.*
- Lee, T.-C., Kashyap, R.L., Chu, C.-N., 1994. Building skeleton models via 3-D medial surface axis thinning algorithms. *CVGIP: Graph. Models Image Process.* 56 (6), 462–478.
- Li, L.H., Yatskar, M., Yin, D., Hsieh, C.-J., Chang, K.-W., 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- Liu, R., Wang, Y., Baba, T., Masumoto, D., 2010. Shape detection from line drawings with local neighborhood structure. *Pattern Recognit.* 43 (5), 1907–1916.
- Long, J., Shelhamer, E., Darrell, T., 2015. Fully convolutional networks for semantic segmentation.
- Mednionogov, A., Kyrki, V., Kälviäinen, H., 2000. Content-based matching of line-drawing images using the Hough transform. *IJDAR* 3, 117–124.
- Mitsubishi UFJ Research & Consulting Co., L., 2019. A survey on projects and issues in Japan's manufacturing industry. [https://www.meti.go.jp/meti\\_lib/report/2020fy/000066.pdf](https://www.meti.go.jp/meti_lib/report/2020fy/000066.pdf).
- Monti, F., Boscaini, D., Masci, J., Rodola, E., Svoboda, J., Bronstein, M.M., 2017. Geometric deep learning on graphs and manifolds using mixture model cnns. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 5115–5124.
- Namata, G., London, B., Getoor, L., Huang, B., Edu, U., 2012. Query-driven active surveying for collective classification. In: *10th International Workshop on Mining and Learning with Graphs*, Vol. 8. p. 1.
- Paliwal, S., Sharma, M., Vig, L., 2021. OSSR-PID: One-shot symbol recognition in P&ID sheets using path sampling and GCN. In: *2021 International Joint Conference on Neural Networks. IJCNN, IEEE*, pp. 1–8.
- Perozzi, B., Al-Rfou, R., Skiena, S., 2014. Deepwalk: Online learning of social representations. In: *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. pp. 701–710.
- Redmon, J., Divvala, S., Girshick, R., Farhadi, A., 2016. You only look once: Unified, real-time object detection. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 779–788.
- Rica, E., Alvarez, S., Moreno-Garcia, C.F., Serratos, F., 2023. Zero-error digitisation and contextualisation of piping and instrumentation diagrams using node classification and sub-graph search. In: *Structural, Syntactic, and Statistical Pattern Recognition: Joint IAPR International Workshops, S+ SSPR 2022, Montreal, QC, Canada, August 26–27, 2022, Proceedings*. Springer, pp. 274–282.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional networks for biomedical image segmentation.
- Sajadfar, N., Ma, Y., 2015. A hybrid cost estimation framework based on feature-oriented data mining approach. *Adv. Eng. Inform.* 29 (3), 633–647.
- Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., Eliassi-Rad, T., 2008. Collective classification in network data. *AI Mag.* 29 (3), 93.
- Shi, J., Malik, J., 2000. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.* 22 (8), 888–905.
- Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*.
- Sousa, P., Fonseca, M.J., 2010. Sketch-based retrieval of drawings using spatial proximity. *J. Vis. Lang. Comput.* 21 (2), 69–80.
- Stutz, D., Hermans, A., Leibe, B., 2018. Superpixels: An evaluation of the state-of-the-art. *Comput. Vis. Image Underst.* 166, 1–27.
- Sulaiman, R., Amran, M.F.M., Abd Majid, N.A., 2012. A study on information extraction method of engineering drawing tables. *Int. J. Comput. Appl.* 50 (16).
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y., 2017. Graph attention networks. *arXiv preprint arXiv:1710.10903*.
- Wang, W., Bao, H., Dong, L., Wei, F., 2021. VLMO: Unified vision-language pre-training with mixture-of-modality-experts. *arXiv preprint arXiv:2111.02358*.
- Wang, J., Sun, K., Cheng, T., Jiang, B., Deng, C., Zhao, Y., Liu, D., Mu, Y., Tan, M., Wang, X., et al., 2020. Deep high-resolution representation learning for visual recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 43 (10), 3349–3364.
- Welling, M., Kipf, T.N., 2016. Semi-supervised classification with graph convolutional networks. In: *J. International Conference on Learning Representations (ICLR 2017)*.
- Xiaojin, Z., Zoubin, G., 2002. Learning from Labeled and Unlabeled Data with Label Propagation. *Tech. Rep., Technical Report CMU-CALD-02-107*, Carnegie Mellon University.
- Xie, L., Lu, Y., Furuhashi, T., Yamakawa, S., Zhang, W., Regmi, A., Kara, L., Shimada, K., 2022. Graph neural network-enabled manufacturing method classification from engineering drawings. *Comput. Ind.* 142, 103697.
- Yang, Z., Cohen, W., Salakhudinov, R., 2016. Revisiting semi-supervised learning with graph embeddings. In: *International Conference on Machine Learning*. PMLR, pp. 40–48.
- Yang, L., Zhuang, J., Fu, H., Wei, X., Zhou, K., Zheng, Y., 2021. Sketchgcn: Semantic sketch segmentation with graph neural networks. *ACM Trans. Graph.* 40 (3), 1–13.
- Zhang, T.Y., Suen, C.Y., 1984. A fast parallel algorithm for thinning digital patterns. *Commun. ACM* 27 (3), 236–239.
- Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J., 2017. Pyramid scene parsing network. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. pp. 2881–2890.
- Zhu, X., 2005. *Semi-Supervised Learning with Graphs*. Carnegie Mellon University.