

Vision Language Meets Transformers

Anonymous ACL submission

Abstract

Recent vision transformers model images as token sequences but rely on patch-based or continuous pixel embeddings. We extend the vision-language paradigm by representing images as textual sequences derived directly from pixels and modeling them with transformer-based language models. Pixel intensities are mapped to Unicode characters and tokenized using Byte Pair Encoding (BPE), after which a BERT-style transformer is trained on the resulting sequences. As a proof of concept, we pre-train on ImageNet-1K and fine-tune on MNIST and Fashion-MNIST. Results show that vision-language-based transformer models can effectively operate on pixel-derived text and benefit from scalable vocabularies, framing images as a discrete and extensible language.

1 Introduction and Related Work

Recent advances in computer vision have enabled high-accuracy performance across tasks such as image classification, detection, and segmentation, driven by deep learning, large-scale datasets, and increased computational resources (LeCun et al., 2015; Goodfellow, 2016; Bengio et al., 2017). Classical vision pipelines treat images as numeric pixel matrices and rely on handcrafted or learned feature extraction (Ram et al., 2013). The introduction of Convolutional Neural Networks (CNNs) marked a paradigm shift by enabling hierarchical feature learning directly from pixel intensities (Li et al., 2021; O’Shea, 2015; Gu et al., 2018; Yamashita et al., 2018; Zeiler and Fergus, 2014), with residual architectures further improving scalability and performance on complex datasets (Targ et al., 2016).

In parallel, natural language processing (NLP) has advanced rapidly with transformer-based architectures (Vaswani, 2017), which model sequences via self-attention and form the basis of modern language models. Pretrained transformers such as BERT (Devlin, 2018) exhibit strong contextual

representations and transferability across tasks. Inspired by these successes, transformers have been adapted to vision by representing images as token sequences. Vision Transformers (ViT) (Dosovitskiy, 2020) and variants such as Convolutional vision Transformers (CvT) (Wu et al., 2021) tokenize images into fixed-size patches, preserving a locality bias from CNNs. More recently, Nguyen et al. (Nguyen et al., 2025) demonstrated that vanilla transformers can operate directly on individual pixels without patch aggregation, challenging the necessity of strong spatial inductive biases. At the intersection of vision and language, multimodal models such as CLIP (Radford et al., 2021) learn joint image–text representations but do not treat image pixels as linguistic units, leaving the question of language-like structure at the pixel level largely unexplored. Moreover, a recent line of work introduced vision-language (Islam et al., 2025), which represents images as textual sequences derived from pixel-level data by mapping pixel intensities to symbolic characters. This enables the use of classical NLP techniques for vision tasks and reveals language-like statistical regularities in pixel-derived sequences. However, the approach remains limited to shallow and lossy architectures, leaving transformer-based language modeling unexplored. **Our Contribution.** We extend the vision-language paradigm to transformer architectures by encoding each pixel of grayscale images as a discrete Unicode symbol, yielding a lossless representation with a base vocabulary of 256 tokens. Images are treated as character sequences, and Byte Pair Encoding (BPE) (Shibata et al., 1999) is applied to learn higher-level visual tokens. A BERT-style transformer is trained directly on these sequences to capture contextual dependencies among pixel-derived tokens.

Unlike patch-based or pixel-token vision transformers, our approach adopts a linguistic abstraction: images are modeled as sequences of discrete

083 symbols rather than continuous embeddings. While
084 ViT-style models are constrained by fixed hidden
085 sizes (e.g., 768) and sequence length, our formula-
086 tion allows representational capacity to scale with
087 vocabulary size via subword tokenization (e.g.,
088 8K–16K tokens).

089 As a proof of concept, we pretrain on ImageNet-
090 1K and fine-tune on grayscale MNIST and Fashion-
091 MNIST, demonstrating that vision-language-based
092 transformer models can effectively operate on pixel-
093 derived text and offering a complementary perspec-
094 tive that frames images as a discrete, extensible
095 language.

096 2 Methodology

097 In this section, we extend the idea of applying
098 vision-language to transformers. Figure 1 illus-
099 trates the whole framework.

100 2.1 Characterization

101 **Image Normalization and Grayscale Projection.**

102 As preprocessing, all images are resized to a fixed
103 resolution of 224×224 pixels and converted to
104 grayscale for color inputs. Let $I \in [0, 255]^{224 \times 224}$
105 denote the resulting image. This transformation col-
106 lapses the RGB channels into a single luminance
107 value, producing a dense matrix of pixel intensi-
108 ties while enforcing a uniform sequence length
109 across the dataset. The grayscale conversion pre-
110 serves essential structural information relevant for
111 downstream visual modeling while reducing in-
112 put dimensionality, yielding a simpler and more
113 tractable corpus for training. As illustrated in Fig-
114 ure 1, the MNIST and Fashion-MNIST datasets
115 are already provided in grayscale and therefore di-
116 rectly processed using 256 intensity levels. For
117 other datasets, images are to be preprocessed into
118 grayscale as a prerequisite for applying the pro-
119 posed framework.

120 **Image-to-Text Characterization.** To transform
121 images into a modality compatible with sequence-
122 based language models, we convert each image
123 into a deterministic sequence of Unicode sym-
124 bols. Each pixel intensity, $p \in \{0, 1, \dots, 255\}$
125 is mapped to a unique character in the Unicode
126 Private Use Area (PUA). We define a mapping:
127 $\{0, \dots, 255\} \rightarrow \{U+E000, \dots, U+E0FF\}$. This
128 results in a closed vocabulary of exactly 256 atomic
129 symbols, each corresponding to a discrete pixel
130 value. The PUA ensures that these symbols do not
131 collide with natural language tokens and remain

132 fully controllable within the framework. Because
133 the mapping operates at the pixel level, considering
134 each pixel as a character, the representation is loss-
135 less for grayscale images and perfectly invertible.

136 **Sequence Linearization.** The grayscale
137 image matrix I is flattened in row-major
138 order into a sequence of length 50,176:
139 $[(I_{1,1}), (I_{1,2}), \dots, (I_{224,224})]$. Thus, each
140 image yields a fixed-length sequence of 50,176
141 Unicode characters, where each character cor-
142 responds to a single pixel. This representation
143 converts the image into a pure text sequence
144 without requiring any learned tokenizer or patch
145 extraction. In contrast to patch-based encodings,
146 this formulation preserves full pixel resolution
147 and enables direct application of sequence models
148 (e.g., masked language modeling or autoregressive
149 objectives) to the raw visual signal.

150 2.2 Tokenization

151 We analyze the subword vocabulary learned by the
152 BPE tokenizer trained over the Unicode-encoded
153 image corpus. Since the input consists of sequences
154 over a fixed 256-symbol alphabet (corresponding
155 to grayscale intensities mapped to Unicode code
156 points), BPE operates by merging frequently co-
157 occurring symbol pairs into composite tokens of
158 increasing length. To characterize the structure
159 of the learned vocabulary, we group all tokens by
160 their character length, where a length-1 token cor-
161 responds to a single Unicode pixel symbol and tokens
162 of length $\ell > 1$ represent ℓ -character substrings
163 formed through one or more merge operations.

164 2.3 Transformer Implementation

165 As illustrated in Figure 1, our approach employs
166 a standard encoder-only transformer architecture,
167 comprising stacked self-attention and feed-forward
168 layers. The proposed vision-language formulation
169 is model-agnostic and can, in principle, be instan-
170 tiated with any transformer encoder capable of pro-
171 cessing token sequences. In this work, we employ
172 BERT as a proof-of-concept to demonstrate the
173 feasibility of modeling pixel-derived linguistic rep-
174 resentations using established NLP architectures.

175 **Pretraining.** We pretrain a BERT-style transformer
176 encoder on the ImageNet-1K dataset, treating each
177 image as a sequence of BPE tokens. The model
178 is trained using the masked language modeling
179 (MLM) objective, where a subset of visual to-
180 kens is randomly masked, and the model is op-
181 timized to predict the original tokens based on

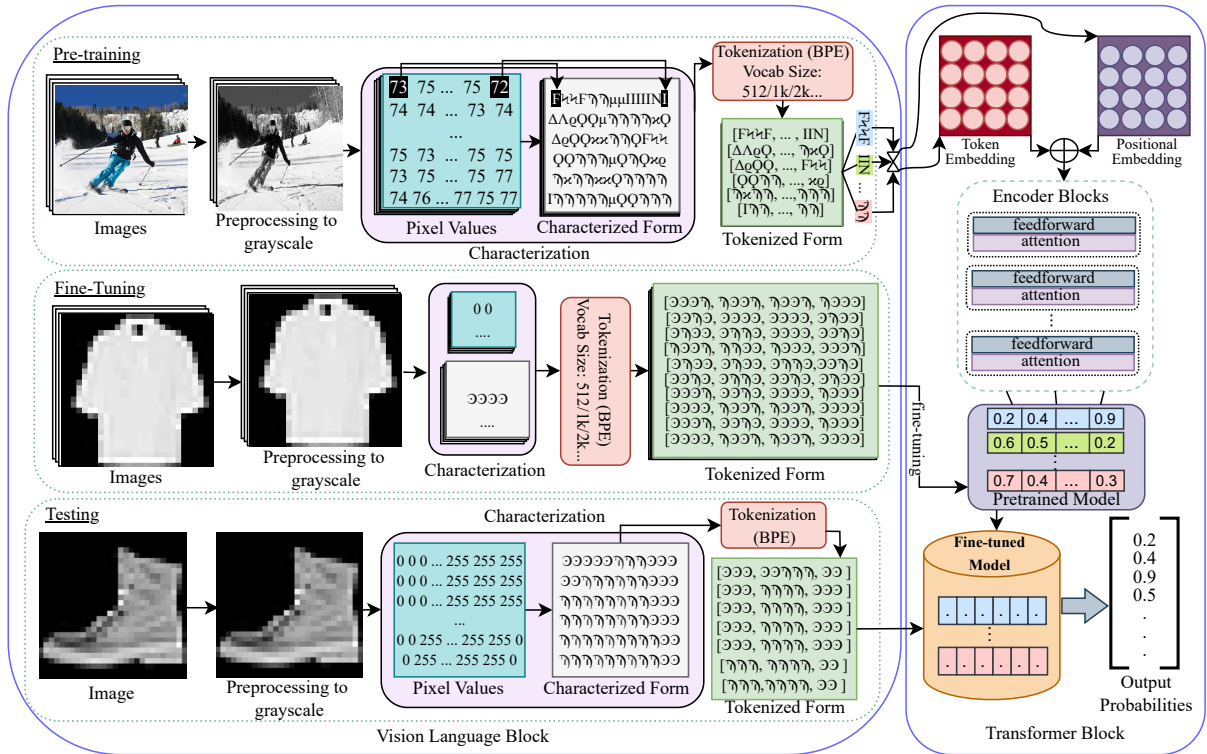


Figure 1: Overview of the vision-language Transformer framework. Images are resized and converted to grayscale, then each pixel intensity is mapped to a unique Unicode character to form a character-level visual text sequence. The resulting sequences are tokenized using Byte Pair Encoding (BPE) and modeled with a transformer trained using a masked language modeling objective, followed by fine-tuning the model.

182 their surrounding context. Through this objective, 183 the transformer learns contextual dependencies 184 among pixel-derived tokens using self-attention, 185 effectively modeling long-range spatial relationships 186 in the image as linguistic context. No label 187 supervision is used during pretraining.

188 **Fine-tuning and Classification.** Following pre- 189 training, the model is fine-tuned on smaller 190 grayscale datasets (MNIST and Fashion-MNIST). 191 Images from these datasets are converted to Uni- 192 code sequences using the same pixel-to-character 193 mapping and tokenized using the pretrained BPE 194 tokenizer. We choose image classification as a 195 downstream task. For classification, a task-specific 196 classification head is added on top of the trans- 197 former, and the representation of the special [CLS] 198 token is used for prediction. Fine-tuning updates 199 all model parameters using labeled data.

200 This implementation demonstrates that a stan- 201 dard transformer encoder, without architectural 202 modification, can effectively model pixel-level 203 vision-language representations. The use of BERT 204 serves to validate the generality of the framework 205 and establishes a baseline for future exploration 206 with alternative transformer architectures.

207 3 Experimental Result Analysis

208 We present a proof-of-concept evaluation of vision- 209 language-based transformer modeling on MNIST 210 and Fashion-MNIST. We compare our lossless ap- 211 proach with the original vision–language lossy 212 baseline (Islam et al., 2025), which reduces im- 213 age resolution via pixel-intensity quantization (e.g., 214 8 or 32 levels) and encodes pairs of vertically adja- 215 cent pixels as a single character. In contrast, we re- 216 present each pixel as a distinct character, preserving 217 full pixel-level granularity and enabling direct ap- 218 plication of standard language modeling techniques 219 without spatial aggregation. For fair comparison, 220 all methods use the same BERT-based transformer 221 architecture with hidden dimension 768, ensuring 222 that observed differences stem from representation 223 choices rather than model capacity.

224 Table 1 reports classification accuracy on 225 MNIST and Fashion-MNIST across different vo- 226 cabulary sizes and pixel quantization levels, re- 227 vealing a strong interaction between pixel granu- 228 larity and vocabulary capacity. For coarse quan- 229 tization (8 shades), high performance is achieved 230 with small vocabularies (256–512), while larger 231 vocabularies consistently degrade accuracy, indi-

Table 1: Classification accuracy (%) on MNIST (M) and Fashion-MNIST (F-M) across vocabulary sizes and pixel quantization levels. The 8- and 32-shade settings represent two vertically adjacent pixels as a single character, while the 256-shade setting treats each individual pixel as a distinct character.

Vocab. Size	8 Shades (Lossy)		32 Shades (Lossy)		256 Shades (Lossless)	
	M	F-M	M	F-M	M	F-M
256	98.50	87.85	–	–	–	–
512	97.69	86.64	99.23	91.77	94.90	83.63
1024	96.70	85.66	99.02	90.89	93.20	83.73
2048	–	–	98.86	90.46	98.00	91.40
4096	–	–	–	–	98.08	92.32
8192	–	–	–	–	98.08	91.96

cating that heavily quantized representations do not benefit from increased lexical capacity. With 32 shades, performance improves substantially, especially on Fashion-MNIST, and intermediate vocabularies (512–1024) perform best, suggesting that finer pixel distinctions require moderate vocabulary capacity. In contrast, the 256-shade setting, corresponding to the pixel-level vision-language formulation, exhibits a markedly different trend. Small vocabularies (512–1024) perform poorly, while larger vocabularies are essential for strong performance. Vocabulary sizes of 2048–8192 yield substantial gains, with the best results achieving state-of-the-art-level accuracy of $\approx 99\%$ on MNIST and $\approx 93\%$ on Fashion-MNIST. Overall, these results highlight a clear trade-off: coarse representations are effective with small vocabularies, whereas fine-grained, pixel-level vision-language representations require larger lexical capacity to capture richer visual token patterns.

Furthermore, Table 2 summarizes the n -gram token length distributions learned by BPE under the 256-shade pixel representation. The number of unigram tokens remains fixed at 256, corresponding to the pixel-level character inventory. As vocabulary size increases, bigram tokens grow sharply (from 240 at 512 to 6543 at 8192), indicating that frequently co-occurring pixel pairs dominate the learned visual lexicon and that larger vocabularies enable finer modeling of local interactions. Higher-order tokens emerge only with sufficient vocabulary capacity. Trigrams and longer n -grams are rare at small vocabularies but become increasingly prevalent at larger sizes, with the 8192 setting learning tokens spanning up to 256 characters. Although sparse, these long tokens capture recurring higher-order visual structures over extended pixel sequences. Compared to the 4096 setting, the 8192

Table 2: Distribution of n -gram token lengths across different vocabulary sizes for 256 shades.

Token Length	512	1024	2048	4096	8192
Unigram	256	256	256	256	256
Bigram	240	732	1653	3382	6543
Trigram	0	2	10	148	929
4-gram	5	18	104	258	295
5-gram	4	4	4	6	17
6-gram	1	1	1	4	8
7-gram	0	0	0	0	4
8-gram	3	5	8	26	103
12-gram	0	0	1	1	2
16-gram	2	3	5	6	15
20-gram	0	0	0	0	1
24-gram	0	0	0	0	2
28-gram	0	0	0	0	1
32-gram	1	2	3	5	6
40-gram	0	0	0	0	1
44-gram	0	0	0	0	1
48-gram	0	0	0	1	1
64-gram	0	1	2	2	4
128-gram	0	0	1	1	2
256-gram	0	0	0	0	1

vocabulary substantially increases both the number and diversity of mid- and long-range tokens. Overall, the distribution exhibits a pronounced long-tail pattern, closely mirroring subword statistics in natural language and explaining the performance gains observed with larger vocabularies in fine-grained pixel-level vision-language modeling.

4 Conclusion

We extend the vision-language paradigm to transformer-based models by representing images as textual sequences derived directly from pixel-level data. By mapping pixel intensities to discrete symbols, applying subword tokenization, and training a transformer, we show that standard NLP architectures can effectively model pixel-derived visual representations. Pretraining on ImageNet-1K followed by fine-tuning on MNIST and Fashion-MNIST provides a proof of concept that pixel-level vision data exhibits learnable language-like structure. While this study focuses on grayscale images, extending the framework to RGB images, with three channels per pixel, would yield longer and more expressive visual sequences, constituting a richer language. Beyond color images, the approach naturally extends to higher resolutions, diverse visual domains, and tasks such as segmentation, detection, generative modeling, and self-supervised learning. Future work may also explore alternative transformer architectures, efficiency-oriented designs, and structured tokenization strategies, highlighting linguistic abstraction as a flexible foundation for vision modeling.

302 Limitations

303 This work is intended as a proof of concept and has
304 limitations. The current experimental evaluation
305 is limited to grayscale datasets and image classi-
306 fication tasks. While the approach generalizes in
307 principle to RGB images and more complex visual
308 domains, additional studies are necessary to assess
309 its effectiveness on higher-dimensional inputs and
310 tasks such as object detection, segmentation, and
311 image generation.

312 Moreover, our implementation relies on a stan-
313 dard BERT-style encoder with fixed hyperparam-
314 eters. We do not explore alternative transformer
315 architectures, tokenization strategies, or optimiza-
316 tion techniques that may yield better efficiency or
317 performance. Addressing these limitations is an
318 important direction for future research.

319 References

320 Yoshua Bengio, Ian Goodfellow, and Aaron Courville.
321 2017. *Deep learning*, volume 1. MIT press Cam-
322 bridge, MA, USA.

323 Jacob Devlin. 2018. Bert: Pre-training of deep bidi-
324 rectional transformers for language understanding.
325 *arXiv preprint arXiv:1810.04805*.

326 Alexey Dosovitskiy. 2020. An image is worth 16x16
327 words: Transformers for image recognition at scale.
328 *arXiv preprint arXiv:2010.11929*.

329 Ian Goodfellow. 2016. Deep learning.

330 Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma,
331 Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing
332 Wang, Gang Wang, Jianfei Cai, and 1 others. 2018.
333 Recent advances in convolutional neural networks.
334 *Pattern recognition*, 77:354–377.

335 Aminul Islam, Md Mustakin Alam, and Shaker Islam.
336 2025. [The birth of vision language](#). In *Proceedings*
337 *of the 33rd ACM International Conference on Multi-*
338 *media*, MM '25, page 12361–12370, New York, NY,
339 USA. Association for Computing Machinery.

340 Yann LeCun, Yoshua Bengio, and Geoffrey Hinton.
341 2015. Deep learning. *nature*, 521(7553):436–444.

342 Zewen Li, Fan Liu, Wenjie Yang, Shouheng Peng, and
343 Jun Zhou. 2021. A survey of convolutional neural net-
344 works: analysis, applications, and prospects. *IEEE*
345 *transactions on neural networks and learning sys-*
346 *tems*, 33(12):6999–7019.

347 Duy Kien Nguyen, Mido Assran, Unnat Jain, Martin R.
348 Oswald, Cees G. M. Snoek, and Xinlei Chen. 2025.
349 [An image is worth more than 16x16 patches: Exploring](#)
350 [transformers on individual pixels](#). In *The Thir-*
351 *teenth International Conference on Learning Repre-*
352 *sentations*.

K O’Shea. 2015. An introduction to convolutional neu- 353
ral networks. *arXiv preprint arXiv:1511.08458*. 354

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 355
Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas- 356
try, Amanda Askell, Pamela Mishkin, Jack Clark, and 357
1 others. 2021. Learning transferable visual models 358
from natural language supervision. In *International* 359
conference on machine learning, pages 8748–8763. 360
PMLR. 361

Idan Ram, Michael Elad, and Israel Cohen. 2013. Im- 362
age processing using smooth ordering of its patches. 363
IEEE transactions on image processing, 22(7):2764– 364
2774. 365

Yusuxke Shibata, Takuya Kida, Shuichi Fukamachi, 366
Masayuki Takeda, Ayumi Shinohara, Takeshi Shino- 367
hara, and Setsuo Arikawa. 1999. Byte pair encoding: 368
A text compression scheme that accelerates pattern 369
matching. 370

Sasha Targ, Diogo Almeida, and Kevin Lyman. 2016. 371
Resnet in resnet: Generalizing residual architectures. 372
arXiv preprint arXiv:1603.08029. 373

A Vaswani. 2017. Attention is all you need. *Advances* 374
in Neural Information Processing Systems. 375

Haiping Wu, Bin Xiao, Noel Codella, Mengchen Liu, 376
Xiyang Dai, Lu Yuan, and Lei Zhang. 2021. Cvt: 377
Introducing convolutions to vision transformers. In 378
Proceedings of the IEEE/CVF international confer- 379
ence on computer vision, pages 22–31. 380

Rikiya Yamashita, Mizuho Nishio, Richard Kinh Gian 381
Do, and Kaori Togashi. 2018. Convolutional neural 382
networks: an overview and application in radiology. 383
Insights into imaging, 9:611–629. 384

Matthew D Zeiler and Rob Fergus. 2014. Visualiz- 385
ing and understanding convolutional networks. In 386
Computer Vision–ECCV 2014: 13th European Con- 387
ference, Zurich, Switzerland, September 6–12, 2014, 388
Proceedings, Part I 13, pages 818–833. Springer. 389