Northumbria Research Link

Citation: Ruta, Dan, Gilbert, Andrew, Aggarwal, Pranav, Marri, Naveen, Kale, Ajinkya, Briggs, Jo, Speed, Chris, Jin, Halin, Faieta, Baldo, Filipkowski, Alex, Lin, Zhe and Collomosse, John (2022) StyleBabel: Artistic Style Tagging and Captioning. In: Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VIII. Lecture Notes in Computer Science, 13668. Springer, Cham, Switzerland, pp. 219-236. ISBN 9783031200731, 9783031200748

Published by: Springer

URL: https://doi.org/10.1007/978-3-031-20074-8_13 <https://doi.org/10.1007/978-3-031-20074-8_13>

This version was downloaded from Northumbria Research Link: https://nrl.northumbria.ac.uk/id/eprint/49791/

Northumbria University has developed Northumbria Research Link (NRL) to enable users to access the University's research output. Copyright © and moral rights for items on NRL are retained by the individual author(s) and/or other copyright owners. Single copies of full items can be reproduced, displayed or performed, and given to third parties in any format or medium for personal research or study, educational, or not-for-profit purposes without prior permission or charge, provided the authors, title and full bibliographic details are given, as well as a hyperlink and/or URL to the original metadata page. The content must not be changed in any way. Full items must not be sold commercially in any format or medium without formal permission of the copyright holder. The full policy is available online: http://nrl.northumbria.ac.uk/policies.html

This document may differ from the final, published version of the research and has been made available online in accordance with publisher policies. To read and/or cite from the published version of the research, please visit the publisher's website (a subscription may be required.)





000	StyleBabel: Artistic Style Tagging and Captioning
001	
002	
003	Anonymous ECCV submission
004	Dapar ID 4076
005	r aper 112 4970
006	
007	
800	Abstract. We present StyleBabel, a unique open access dataset of natural lan-
009	guage captions and free-form lags describing the artistic style of over 155K digital artworks, collected via a novel participatory method from experts studying at spe-
010	cialist art and design schools. StyleBabel was collected via an iterative method,
011	inspired by 'Grounded Theory': a qualitative approach that enables annotation
012	while co-evolving a shared language for fine-grained artistic style attribute de-
013	scription. We demonstrate several downstream tasks for StyleBabel, adapting the
014	embeddings for: 1) free-form tag generation: 2) natural language description of
015	artistic style; 3) fine-grained text search of style. To do so, we extend ALADIN
016	with recent advances in Visual Transformer (ViT) and cross-modal representation
017	learning, achieving a state of the art accuracy in fine-grained style retrieval.
018	Konwords, Detects and avaluation. Image and video retrieval Vision 1 lan
019	guage. Vision applications and systems
020	gauge, vision appreadons and systems
021	
022	1 Introduction
023	
024	Artistic style is the distinctive appearance of an artwork; i.e. how an artist has depicted
025	their subject matter [14]. Describing the artistic style of digital artwork is an open chal-
026	in the style domain. However, sutomated style description has notantial applications in
027	summarization analytics and accessibility. For the first time, this paper shows that a
028	set of descriptive tags, or even complete caption sentences, may be automatically gen-
029	erated to describe the fine-grained artistic style of an image – distinct from its content
031	[28,10,55], or the emotions it evokes [2]. Our core contribution enabling this is StyleBa-
032	bel, a novel dataset of fine-grained [54,51,32] annotations describing the artistic style
033	of ~ 135 K digital artworks ¹ , collected from expert participants via a novel participatory
034	method that forms a further contribution of this paper. Specifically, our novel contribu-
00+	tions are:

1. StyleBabel Dataset. We present a new dataset of over ~135K {image, tags, natural language (NL) caption } pairs for digital artwork images annotated by a combination of crowd-sourcing and 48 domain experts drawn from design and art schools. StyleBa-bel is the first dataset containing all three data types for every data sample. Furthermore, the images contained have vastly greater style variation than existing datasets [42,16] which predominantly focus on small subsets of fine art, sometimes further limited to only European or Asian [3,27,53]. Not only is StyleBabel's domain more diverse, but our an-notations also differ. StyleBabel focuses on the visual appearance of images (which can include stroke/colouring type, lighting, shading, patterns, shapes, composition, medium,

 $^{-1}$ The dataset will be released for open access (CC-BY 4.0).



Fig. 1. Results *generated* from models trained on StyleBabel to describe style using free-form tags (left) and captions (right). The top row of tags contains tags generated by our models trained with StyleBabel. The middle row contains tags from a commercial image tagging Google product, and the bottom from a similar commercial product from Microsoft. We show the importance of the style domain-specific information within StyleBabel, for generating relevant style captions.

layout, theme etc.), we avoid external, high level information such as artists, time peri ods, surrounding meaning, emotions evoked, provenance facts, context or content. This
 enables new avenues for research not possible before, some of which we explore in this
 paper.

We do not seek to align our work to a formal ontology or definition of style (some-thing heavily debated even by academics [14,35]). Instead, we explicitly build, evolve, and rely on the emergent structure of style information from the collective, harmonized experience of expert collaborators from art and design universities during our data col-lection process. In previous work [17], the working definition of style has been the similarity of gram matrices at specific layers in CNN backbones such as VGG [45]. Dur-ing our data pre-processing/initialization, ahead of style annotation, we similarly use a working definition as similarity in the ALADIN [41] style model's embedding space, previously shown to accurately represent a variety of artistic styles in a metric space.

During annotation, images are grouped into $\sim 6K$ moodboards (grids of style-consistent images). Each moodboard is annotated by a group of experts, both with tags and freeform captions, to yield a description of visual style. The vocabulary used for both annotation forms is unconstrained. The moodboard annotations are cross-validated as part of the collection process and refined further via the crowd to obtain individual, image-level fine-grained annotations.

2. Grounded Annotation Methodology. We present a new data annotation method-ology inspired by Grounded Theory (GT) [44,6], a qualitative research method used in the humanities and social sciences. In GT, participant groups engage in an unconstrained data clustering exercise while simultaneously evolving a shared vocabulary for describ-ing those clusters. Working with these disciplines, we adapted this process to evolve a shared vocabulary for annotation while annotating groups ('moodboards') of images with similar artistic styles. The moodboards are obtained by clustering artworks within the ALADIN fine-grained embedding for style similarity [41]. Our iterative methodol-ogy comprises individual and participatory group stages and a validation stage.

3. Artistic Retrieval and Description. We incorporate a Visual Transformer [46]
 into a fine-grained visual style representation architecture [41], (ALADIN-ViT). ALADIN-

ViT provides state of the art performance at fine-grained style similarity search. We train models for several cross-modal tasks using ALADIN-ViT and StyleBabel annotations. Using CLIP [36], we train with StyleBabel to generate free-form tags describing the artistic style, generalizing to unseen styles. We show that we may apply these tags for text based style search. We similarly demonstrate the synthesis of descriptive natural language captions for digital art.

2 Related Work

Representation learning for visual style has focused primarily on neural style transfer(NST) and style classification.

Style Transfer. Classical approaches learned patch-based representation of style by analogy from paired data. Gatys et al. [17] enabled NST by extracting disentangled representations separating content and style from unpaired data, using a Grammian computed across layers of a pre-trained VGG-19 [45] model. Similarly, feed-forward networks used the Grammian to train fast encoder-decoder models for NST specialized to pre-trained styles [47,23]. Extensions to multi-scale [50] and video [40] NSTS were later presented. To relax the constraint to pre-trained styles, feature based NST was proposed using Adaptive Instance Normalization (AdaIN) [21,18]. This approach was further generalized by the whitening and coloring transform (WCT), which matched feature covariances [29]. Recently unsupervised style transfer was enabled via MUNIT [22] and swapping autoencoder [34]. The latent style codes of encoder-decoder networks for NST were recently adapted for style-based visual representation learning, using weak supervision to learn a metric embedding for fine-grained similarity using the ALADIN model [41].

Style analytics via classification has been explored for digital artwork [52], smaller fine-art collections [57,24], product designs [4], and to identify both painters [7] and genres [43]. Style-based visual search has also been proposed for coarse-grained style using triplet [9] and constrastive training for fine-grained style via ALADIN [41].

Style datasets. No prior dataset of fine-grained artistic style description exists. How-ever, several annotated datasets of artwork have been produced. Behance Artistic Media - BAM [52] comprises 2M diverse digital artworks from the Behance.net platform. with 7 coarse-grained style and 4 emotion tags and no descriptions. **Omniart** (432K im-ages) and Wikiart (81K) are datasets of fine art with associated metadata but no descrip-tions or tags. The **SemArt** dataset [16] focuses on very high level contextual semantic information, rarely containing style (visual appearance) information, and narrowly fo-cuses solely on 8th-19th century European fine art. The BAM-FG (BAM Fine-grained) dataset comprises 2.62M images grouped into 310K style-consistent groupings but no descriptive text or labels. Our work also uses Behance.net, to provide expert style tags and natural language descriptions over 135K images. The AVA dataset [5] studies aesthetics information in *photographic* images only, and the follow-up **AVA-captions** dataset [19] adds captions for these, sourced from noisy internet comments sections. Recently, ArtEmis [2] released non-expert annotations capturing the emotions felt by viewers of fine art in WikiArt. Our proposed StyleBabel dataset is aligned to this contem-porary work in that it also seeks to ascribe text to visual art. However, our focus is also on digital, not just fine artwork. We also differ in that our annotation is led by expert students in specialized art and design schools, not exclusively by non-expert crowd-workers. Notably, our annotations focus on the style alone, deliberately *avoiding* the

description of the subject matter or the emotions that matter evokes. ArtEmis instead describes *exclusively the emotions evoked* by both the style and content of the artwork. both of which feature in the descriptions. To recap. StyleBabel is unique in providing tags and textual descriptions of the artistic style, doing so at a large scale and for a wider variety of styles than existing datasets, with labels sourced from a large, diverse group of experts across multiple areas of art.

Image captioning and visual question answering methods [49] initially learned LSTM language models, leveraging semantic image embeddings e.g. via ResNet/ImageNet. Later image captioning work [28,10,55,56] made use of object detection: regions of in-terest (ROI). Their relationships vielded improved semantics captioning models, though often due to the bias of co-present context that hinted at the image narrative. This is incompatible with our domain of artistic style, where this localization bias is not some-thing we can use. Recently, there has been an influx of research combining the visual and text domains [36,37,38]. Established methods [13,25], have primarily functioned by generating captions from templates. Though this approach benefits from generat-ing grammatically correct captions, its inflexibility has led to diminished interest in its application, despite recent implementation by OpenAI's CLIP [36]. Here, captioning capabilities generate tags and insert them into templates, using two encoders – for text and image modalities. This model then performs cross-modal training via contrastive loss. More recently, Attention-on-attention [20] can generate state-of-the-art quality captions from an image embedding by filtering out irrelevant attention results. VirTex [11] recently demonstrates that caption annotations are more efficient for representa-tion learning of images, with better or comparable representation quality on ImageNet despite much less training data.

Grounded Theory (GT) [44] is a qualitative method used in the humanities and social sciences to codify data - i.e. to identify and name apparent or emergent pat-terns across different data sources and types. Through interpretation, participants col-laboratively agree on an initial set of summative 'open' codes that are then iteratively combined into larger groups, until eventually, participants arrive, by consensus, at the common themes in the data [6]. We adopt this approach to collecting expert annotation to describe the artistic style of clusters of digital artwork. Free-form textual input from various participants can vary in writing style, creating a very noisy dataset. GT mitigates this by guiding participants to align their responses to a consistent format.

StyleBabel Dataset

StyleBabel is a new dataset for cross-modal representation learning. It comprises 135k digital artwork images from the public creative portfolio website Behance.net (in turn, available via the BAM dataset). Each image is annotated with a set of keyword tags and natural language descriptions 'captions' describing its fine-grained artistic style -the distinctive appearance of the image - in the English language. We focus mainly on attributes that we can visually depict, rather than more high level and abstract concepts such as emotions [2]. StyleBabel enables the training of models for style retrieval and generates a textual description of fine-grained style within an image: automated natural language style description and tagging (e.g. style2text). We train state of the art proof of concept models for these tasks using our dataset in Sec. 5.

180 3.1 Study Context

We determined via early initial trials with crowd annotation platforms (AMT) that the quality, coarseness, and diversity of data generated by non-experts is indequate for style description tasks. We collaborated with graduate schools specializing in digital art and design to address this. Together with academic experts at these schools, we designed a novel multi-staged participatory method to enable novel style vocabulary gathering. tagging, and caption generation, recruiting 48 expert staff and student participants. We particularly sought (but did not make a prerequisite) participants familiar with Behance. The final cohort comprised a representative balance of gender and ethnic background of both graduate and undergraduate cohorts - anecdotally, the gender split was equal. Expertise was at that of a final year undergraduate student, with more senior faculty members participating in the discussions. Their program's specialisms were primarily communication design (graphics, illustration), industrial design, fashion, and animation.

We executed the group exercise online using a *collaborative whiteboard* interactive platform², that provides capabilities for multiple users to move and annotate components freely in real time. We built 1300 unique pages of moodboards (Fig.2) to allow partic-ipants to move sticky notes with related tags naturally in a collaborative process. The process simulated in-person group collaboration, enabling the formation of tag clusters and aligning to processes of GT codification. The combined use of Miro and Zoom supported real-time spatial organization of information and associated discussion. Workers were paid significantly higher than the national minimum, with a total dataset cost of approximately \$160k, which we freely contribute as CC-BY 4.0.

3.2 StyleBabel Grounded Annotation

StyleBabel *does not* aim to develop an ontology to categorize style, with agreement on a diverse ontology eluding art practitioners [14]. Yet, consistency of language is essential for learning of effective representations. Therefore we propose an annotation methodology that enables annotations at scale (multiple participants) and encourages co-evolution of a harmonized natural language to describe the style.

Our annotation process instead is inspired by Grounded Theory (GT) [44,6]; a qual-itative method often used for data analysis in the humanities and social sciences. A systematic research process to 'codify' empirical data, identify themes from the data, and associate data with those themes. This process is distinct from fitting (or 'annotating') data to pre-existing categories. GT is an iterative process in which participants co-evolve a language to describe the data as they work on clustering and labeling it with that shared language. Concretely, GT often begins with a discussion around a subset of the data during which clusters are formed. Data is moved freely between clusters during the debate, from which a shared understanding and, ultimately, a shared termi-nology evolves for describing those clusters. With further data, this language identifies and names patterns apparent or emergent in the data.

Grounded Annotation Process We propose a multi-stage process for compiling the
 StyleBabel dataset comprised of initial individual and subsequent group sessions and
 a final individual stage. Each batch of these sessions was estimated to take around 10
 hours, run over a week, and run four sets over four weeks.

² https://miro.com/



Fig. 2. The StyleBabel Grounded Annotation process. Moodboards (examples, left) are annotated via an iterative process (right) that encourages groups of participants to evolve and converge to a shared language for describing style. See subsec. **3.2** for stage descriptions.

Experts annotate images in small clusters (referred to as image 'moodboards'). Moodboards are obtained by automatically clustering artworks within a fine-grained style embedding [41]. Our annotation process thus pre-determines the clusters for expert annotation. Still, it encourages expert groups to evolve a harmonized language during the iterative annotation process (as in GT) to improve data consistency. We refer to this process as 'grounded annotation'.

Pre-processing - Moodboard generation (Automated)

We downloaded an initial dataset of 150 million digital artwork (static image assets) from Behance net. We encode all images into a metric search embedding [41], that we then clustered into 6.5K clusters, with L_2 distance to identify the 25 nearest image neighbors to each cluster center. The images from each cluster were arranged in a 5×5 grid for presentation as a moodboard. Thus, we start the annotation process using 6,500 moodboards (162.5K images) of 6.500 different fine-grained styles.³. The extremely high data density from this internet-scale data corpus ensures that the small clusters formed are very stylistically consistent. Fig. 2 (left) shows examples of moodboards.

Stage 1a - Contrastive Annotation (Individual)

Participants were individually presented with a pair of 5x5 moodboards and asked to generate a list of textual tags ('style attributes') that occur in one moodboard, but not another, and a list of style attributes which are shared by both. The moodboards were sampled such that they were close neighbors within the ALADIN style embedding. In the annotation, comparative language (e.g. 'X is brighter than Y') was not permitted to encourage standalone descriptions (e.g. 'bright' was allowed). This paired approach encouraged the suggestion of fine-grained style attributes, supporting participants to suggest characteristics that may otherwise not have been considered when looking at individual styles. Multiple participants annotated each moodboard in this way to produce a rich initial set of attributes.

²⁶² Stage 1b - Harmonization and Refinement (Group)

A collaborative workspace was synthesized within Miro, in which 5 moodboards and their associated style tags from Stage 1a are displayed (as 'sticky notes' below each moodboard).

After an initial briefing and group discussion, each group considered moodboards collectively, one moodboard at a time. All participants were asked to add new tags to ³ We redacted a minimal number of adult-themed images due to ethical considerations

270		一派机器的 过度			李鲁朝		270
271				A	A NS		271
272				N. 180	1. 10 100		272
273				MA V	j2 74 33		273
274		AND A BUTAN AND		A B	應, 余 奉		274
275				影查	AMERICA		275
276	Changed (Re-	line, both have the san	ne white back-	b+w bla	ack and white	, pen and in	ık, fan- ²⁷⁶
277	moved or	grounds with very litt	le tonal varia-	tasy, int	ricate, white s	pace, centra	al com-277
278	added tags in Stage 1b)	tions/ just block colours	s, no difference	position	i, linear, illust	ration	278
279	Stage ID)	constructed by lines n	o colour hand				279
280		drawing linear, pen and	ink				280
201	Final tags after	ink work, sketc	ches, bright,	graphic	, drawing,	black and	white, 281
202	Stage 1b clean-	black and white, white	background,	pen and	ink, fantasy	, intricate,	white 282
200	ing	clean, monotonal, dra	wing, simple,	space,	central con	nposition,	linear,
285		<u>illustration</u> , <u>linear</u> , <u>pen a</u>	and ink	illustrat	<u>101</u>		285
286							286
287	Fig. 3. Before and	d after harmonization (Stag	ge 1b), showing	the bene	fits of the GT a	approach. Ty	wo 287
288	moodboards rece	ive linguistically different	tag sets in Stage	1a (indi	vidual) but are	e conformed	to 288
289	and the tag sets d	ifter but draw upon a share	pant group in Sta ed vocabulary (ur	ige 10 (g	roup). The mo	ouboard styl	$\frac{100}{200}$
290	being shown in th	ne same session to workers	s. Tags removed	are in re	d, and additior	ns are in gree	en. 290
291	C		C			0	291
292	the pre-populate	ed list of tags that we ha	d already gathe	ered from	m Stage 1a (t	he individu	ial 292
293	task), modify th	ne language used, or ren	move any tags	they ag	reed were no	t appropria	ite. 293
294	Each moodboar	d was considered 'finisl	hed' when no n	nore cha	anges to the t	ags list cou	ld 294
295	be readily deter	mined (generally withi	in 1 minute). W	orkers	spent at mos	t 2h 15m p	er 295
296	session, includi	ing all breaks. At least (one facilitator v	vas aiw	ays present t	nrougnout	to 296
297	this part of the s	study via the Miro platfo	orm Fig 3 prov	ides an	illustrative ex	ample of la	ng 297 an-
298	guage harmoni	zation. Two Stage 1a m	oodboards with	n initial	ly very differ	ent langua	298 ge
299	but similar style	e resulted in similar attr	ibutes post-Sta	ge 1b.	5 5	e .	299
300	Stage 2 - Valida	ation and Description (I	Individual)				300
301	The participant	s completed the final sta	age individually	v. Each	participant w	vas present	ed
302	with a random c	collection of 5 moodboar	rds and a set of t	tags gen	erated during	the previo	us 302
303	sessions for just	t one of these moodboar	ds. Participants	s engage	ed in ESP-lik	e game [1]	to 303
205	identify which	moodboard the tags had	d been generate	ed to de	scribe, to ver	rify accura	cy. 304
305	Further, we then	asked them to create na	atural language	caption	s, using as ma	any present	ed 305
307	tags as possible	. We stressed to include	as many tags a	is possit	ble as by now	, they had a	all 307
308	Store 3 Seale	y cleaned and relined.	dividual)				308
309	Stage 5 - Scale	-up una Kejinemeni (Ind					309
310	Following thes	e sessions with our sul	bject experts, v	ve ran a	a different ta	isk with no	on- 310
311	experts (worke	rs on Amazon Mechan	ical Turk (AM	(1)). We	e asked at an	i <i>image lev</i>	vel 311
312	considered not	appropriate Though the	e moodboards	resenta	f to uscard a d to these no	any tags (n m-expert p	cy ar- 312
313	ticipants are sty	de-coherent. there was s	still variation in	the im	ages, meanin	g that certa	$\frac{1}{10}$ $\frac{313}{10}$
314	tags apply to m	lost but not all of the in	nages depicted.	This st	ep helped us	to refine the	he ³¹⁴



Fig. 4. Virtual environment for collaborative Stage 1b. Style attributes (as 'sticky notes') are added, modified and removed from each moodboard to harmonize and filter language.

final tags to individual images further. Unlike novel vocabulary generation, this verification step can be readily performed by non-experts, as detecting invalid tags is a much easier task than novel tag generation. We collected 3 responses per image and performed majority consensus to determine the final tags to use at the per-image level.

We additionally performed a large-scale crowd annotation exercise to individualize cluster-level captions to individual images. Trained workers were presented with individual images, its tags, and the moodboard caption and were asked to compose (potentially many) natural language captions using the tags and caption, ensuring the full set of tags were incorporated across those sentences. A constant set of workers were trained with feedback, for several months, together with a Quality Control (QC) process to ensure high quality annotation. The QC process included grammatical correctness checking, and rejection of any description of content or emotion the descriptive captions.

Language processing Aside from the crowd data filtering, we cleaned the tags emerging from Stage 1b through several steps, including removing duplicates, filtering out invalid data or tags with more than 3 words, singularization, lemmatization, and manual spell checking for every tag. The spell checking step was carried out in 3 passes. The accuracy of the validation step in Stage 2 was found to be 90%.

The final StyleBabel dataset contains 135k images with an average of 12.8 tags per image, over 6k style groups (of the 6,500 initially sampled, with 6k completed by workers in the available time). The tags dictionary contains 3,151 unique tags, and the captions contain 5,475 unique words. Prepositions, determiners, and conjunctions were filtered out. Cardinals (eg. 80s) are kept. 9 cardinals, 2098 nouns, 500 verbs, 55 adverbs, and 482 adjectives are captured. Fig. 5 visualizes a word cloud from the 250 most common style attributes in StyleBabel, and Tbl. 1 shows the richness of the tags and captions.

4 Visual Embedding (ALADIN-ViT)

ALADIN is a two branch encoder-decoder network that seeks to disentangle image content and style. It works by pooling Adaptive Instance Normalization (AdaIN) statistics
 across multiple layers in the style encoder to produce an embedding for fine-grained
 style similarity using a VGG-19 convolutional backbone for the style encoder. In our

	And the second s
Fig. 5. V	isualizing the top 250 tags captured within the StyleBabel annotation over 135K images
Image	
Tags	dim, concept, action, fan- dim, concept, action, fan- tasy, powerful, digital, oil, painting, drawing, bright, colors, draw- photography, animated, pro- artistic, melancholic, ing, child, stroke, drawing, sketch, figure totype, masculine, detailed, pleasing analog professional, lighting pleasing busy, clear, illustra- tion, festive, blue commercial, white scamp, stroke, product pencil, rough, thin, iso metric
Caption	Fantasy themed digital Portrait oil painting Digital bright fantasy Analog experimental illustration featuring an of a female charac-anthropomorphism sketches with thin pen animated male character, ter featuring abstract cartoon illustration cil strokes and lines dim highlighting and a shapes and psychedelic created with soft. The isometric drawing hazy, dark and cluttered patterns against a dark diffused blended hues, expresses commercial background. The illustra-background. The artistic brush strokes, lines, product development. tion highlights the powerful artwork is melancholic and geometric forms masculine character with and using thin repetitive in neutral and cool strokes and shades.
Table 1. collected	Excerpt of StyleBabel dataset. Four images and the corresponding tags and captions i via our Grounded Annotation.
later exp domain we refe lowing weakly Having AdaIN	periments, we require to use a Visual Transformer [12] (ViT) model for the visior. To achieve this, we adapt ALADIN to use ViT as the style encoder backbone r to this as ALADIN-ViT (Fig. 6). We retrain ALADIN-ViT on BAM-FG fol the same training method [41], i.e. using both the reconstruction loss and the supervised contrastive loss and using the implicit style grouping in BAM-FG swapped the style encoder for a transformer, it is no longer possible to sample statistics from feature maps in the encoder. However, we retain the use of the

We achieve the state of the art fine grained style retrieval accuracy on the BAM-FG test partition (Tbl. 2) at 64.48 Top-1, beating not only ALADIN (58.98) but also their fused variant (62.18), which incorporates ResNet embeddings into a concatenated embedding. We, in part, attribute the gains in accuracy to the larger receptive input size (in the pixel space) of earlier layers in the Transformer model, compared to early layers in CNNs. Given that style is a global attribute of an image, this greatly benefits our domain as more weights are trained on more global information.



Fig. 6. ViT-ALADIN architecture used for StyleBabel experiments; as per ALADIN [41] but swapping the style encoder for ViT [12] (change in green) and retrained end-to-end on BAM-FG.

Data	Model	IR-1
BAM-FG-C ₅	ALADIN	58.98
BAM-FG-C ₅	ResNet	45.22
BAM-FG-C ₅	ALADIN (Fused)	62.18
BAM-FG-C ₅	ViT	57.91
BAM-FG-C ₅	ALADIN-ViT	64.48

Table 2. Fine grained style retrieval on the BAM-FG dataset [41] of proposed method ALADIN-ViT, compared to previous methods

StyleBabel Experimental Setup

Data Partitions. We define train/validation/test partitions within StyleBabel for our experiments as follows. Splits were separated on a moodboard basis, to avoid over-lap. The training split has 133k images in 5.974 groups with 3.167 unique tags at an average of 13.05 tags per image. The validation and test splits contain 1k unique images for each validation and test, with 1.256/1.570/10.86 and 1.263/1.636/10.96 unique tags/groups/average tags per image. Captions have an average of 2.4 sentences, with an average of 19.3 words, from 6119 unique words. 1000 samples were extracted for each of the test/validation splits.

Implementation. The models were trained with a maximum batch size of 11k on a 12GB GTX Titan, with a learning rate of 0.003, Adam optimizer, and weight decay of 1e-6. Logit accumulation [41] was employed to reach the maximum batch size possible in the GPU VRAM capacity. We trained all models to convergence. The learning rate was decayed using cosine annealing, as per SimCLR [8]. For the VirTex captioning experiments, a batch size of 105 was used, on a V100, with a learning rate of 2e-4.

Experiments and Discussion

We illustrate the potential of our StyleBabel dataset for three cross-modal learning tasks: **1. Style Auto-tagging:** (style2text) Using StyleBabel tags, we train a CLIP [36] model to learn a cross-modal embedding between image and text embeddings. We generate several tags for unseen StyleBabel images and explore the generated tags' accuracy and the model's ability to generalize.

2. Style Description: (style2text) We similarly make use of the natural language cap-tions collected in StyleBabel, and showcase the generation of natural language captions describing the style of images.



5 tags shown for each image.

3. Text Based Style Retrieval: (text2style) We explore the efficacy of our tag generation for text based style search.

6.1 Style Auto-tagging (style2text)

Recent literature in image captioning has transitioned to making use of object detectors in their model pipelines. This makes sense in semantics, as such features are most often localized to a subset of the image. Style, however, is typically a global attribute of an image, and object detectors are not compatible. Style is more abstract and seldom localized to any specific region of an image. We use the CLIP [36] training methodology to learn a joint embedding space between the publicly available CLIP text encoder and our new vision transformer (ALADIN-ViT). CLIP is traditionally formed of two transformers, the first for text encoding and the second for image encoding. Two MLP heads are trained together through contrastive loss to learn a joint text/image embedding, adding invariance to the modality. We freeze both pre-trained transformers and train the two MLP layers (ReLU separated fully connected layers) to project their embeddings to the shared space. We follow CLIP, employing contrastive loss to drive learning, with the same training set-up.

When using the model for inference, we pass the entire dictionary of available tags through the text encoder and multi-modal MLP head to generate text embeddings. Next, we infer the image embedding using the image encoder and multi-modal MLP head, and calculate similarity logits/scores between the image and each of the text embeddings. We evaluate accuracy via WordNet similarity [15] using the *nltk*⁴ library to first compute synsets for the N ground truth tags for an image. Next, we sort the tags in the dictionary by the logit scores following the embedding inference similarity. We select the top N"retrieved" tags, and for each, and calculate the WordNet similarity to each ground truth tag. The similarity ranges from 0 to 1, where 1 represents identical tags. We save the top

⁴ https://www.nltk.org/

495	Data	Model	WordNet score	495
196	CLIP Webscale	CLIP [36] baseline	0.168	496
407	StyleBabel-mturk	ALADIN-ViT	0.164	400
497	StyleBabel (coarse)	CLIP [36]	0.187	497
498	StyleBabel (coarse)	ALADIN-ViT	0.225	498
499	StyleBabel (FG)	CLIP	0.215	499
500	StyleBabel (FG)	ALADIN-ViT	0.352	500

Table 3. Ablation experiments for tag generation under the CLIP training setting. We show the ben efit of annotating StyleBabel via the proposed GT annotation versus non-expert crowd-sourcing
 (StyleBabel-mturk). We further demonstrate the benefits of tag annotations individualised to the
 image-level (FG), compared to cluster-level (coarse).

value (the similarity score for the closest/most relevant ground truth tag) and repeat it
for each test image. These values are averaged to form our *WordNet score*. We use this
instead of precision/recall, as multiple tags in the dataset can represent very similar (but
not identical) concepts. A soft score better encapsulates the perceived accuracy of the
tags.

We experiment with training CLIP on variants of StyleBabel, presenting results in Tbl.3. In particular, we quantify the difference in quality between collecting annotation via non-expert crowd annotation (StyleBabel-mturk) and gathering expert annotations using our GA process (StyleBabel/ALADIN-ViT). We also show the value of the final stage, where we refine tags to the image-level (FG) rather than moodboard-level (coarse). We train the MLP heads atop the CLIP image encoder embeddings (the 'CLIP' model) and atop embeddings from our ALADIN-ViT model (the 'ALADIN-ViT' model). The former is not based on the ALADIN-ViT style embedding and underperforms by 40%. The best performing model is the proposed ALADIN-ViT trained via StyleBabel data collected using GA on FG labels, with a WordNet score of **0.352**, double the CLIP [36] baseline. Fig. 7 shows some examples of tags generated for various images, using the ALADIN-VIT based model trained under the CLIP method with StyleBabel (FG).



MSa woman with a paint- a close up of a remote a kite that is hanging a room with a lot of a person riding a surf-COCO ing of a face on it with a remote in the air windows and a clock board in the water StyleBabel digital illustration of product photography of digital illustration architectural photog-abstract watercolor CL+IL a fictional character business cards with text featuring animated raphy of a modern in-painting using paint dark color and logo in soft light-characters and typog-terior design in bright brush strokes using and tones against a black ing against a white back-raphy using bright lighting using a neu-shading effects tral color palette background ground colors

Fig. 8. Examples of natural language captions generated for various art styles; Generated captions are compared from VirTex models trained on MS-COCO and Stylebabel CL+IL, showing the benefit of the style information present in the dataset

We run a user study on AMT to verify the correctness of the tags generated, presenting 1000 randomly selected test split images alongside the top tags generated for each.

For each image/tags pair, 3 workers are asked to indicate tags that don't fit the image. We score tags as correct if all 3 workers agree they belong. The absolute accuracy of this study is 89.86%, indicating high tag generation accuracy.

Finally, we explore the model's generalization to new styles by evaluating the aver-age WordNet score of images from the test split. In Fig.9, we group the data samples into 10 bins of distances from their respective style cluster centroid, in the style embedding space. As before, we compute the WordNet score of tags generated using our model and compare it to the baseline CLIP model. Though the quality of the CLIP model is constant as samples get further from the training data, the quality of our model is sig-nificantly higher for the majority of the data split. At worst, our model performs similar to CLIP and slightly worse for the 5 most extreme samples in the test split. But for the majority of the test data, our model considerably outperforms CLIP.

Data	Model	Bleu-1	Bleu-2	Bleu-3	Bleu-4	METEOR	Rouge-L	CIDEr
MS-COCO baseline	VirTex	0.162	0.053	0.016	0.005	0.037	0.145	0.022
StyleBabel (CL)	VirTex	0.127	0.049	0.022	0.010	0.054	0.135	0.076
StyleBabel (IL)	VirTex	0.331	0.187	0.113	0.071	0.129	0.288	0.350
StyleBabel (CL+IL)	VirTex	0.335	0.189	0.118	0.078	0.131	0.288	0.372
StyleBabel (CL+IL)	ResNet LSTM	0.087	0.021	0.008	0.002	0.033	0.080	0.017
StyleBabel (CL+IL)	ALADIN-ViT LSTM	0.094	0.030	0.013	0.006	0.042	0.089	0.034
VirTex	Artemis	0.185	0.083	0.041	0.023	0.081	0.182	0.146
VirTex	Artemis (SB)	0 120	0.031	0.013	0.005	0.034	0.108	0.029

Table 4. (top) VirTex caption generation metrics on 1k holdout StyleBabel test data. CL represents cluster-level, and IL represents image-level. We also run CL+IL, where we fine-tune a clusterlevel model with image-level data labels. (middle) Results with a baseline LSTM model trained over either ResNet or ALADIN image embeddings (bottom) Additional experiments showing performance of a VirTex model trained on an existing dataset (Artemis), and also evaluated on the StyleBabel (SB) test set.





Fig. 9. (left) Generalization experiment for tag generation using baseline CLIP [36] and our StyleBabel trained model using ALADIN-ViT. The test set was sorted by distance in the style embedding space to closest training cluster. Their WordNet scores were binned into 10 quantized distance bands. The numbers atop the bars indicate the number of samples in the corresponding bin. (right) Top 5 style retrieval using textual tags. Using tags generated by ALADIN-ViT/CLIP over the StyleBabel test partition, we perform a keyword based search for artistic style.

585 6.2 Style Description (style2text)

586 We explore the feasibility of using the StyleBabel dataset for generating natural lan-587 guage captions. We conduct our experiments using a fusion of the VirTex [11] backbone 588 for the visual representation learning of image/caption pairs, and Attention on Atten-589 tion (AoA) [20] for the caption decoding. VirTex encodes images without using scene graphs, therefore avoiding issues related to style not being localized in an image. VirTex 590 replaces the Faster-RCNN [39] component in AoA to generate feature maps. We use 591 the pre-trained VirTex on COCO [31], and fine-tune the entire setup, end-to-end on final 592 StyleBabel captions. As per standard practice, during data pre-processing, we remove 593 words with only a single occurrence in the dataset. Removing 45.07% of unique words 594 from the total vocabulary, or 0.22% of all the words in the dataset. We test the caption 595 generation on the StyleBabel test data across Bleu [33], METEOR [26], Rouge [30]. 596 and CIDEr [48], as shown in Tab 4. See the supplementary material for further analysis.

585 586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616 617

618

619

620

621

622

623

624

625

626

627

628

629

We would like to stress that these metrics are not comparable with values from these metrics on standard literature, as we are solving a new task. In literature, these metrics are used for semantic, localized features in images, whereas our task is to generate captions for global, style features of an image. We include an MS-COCO baseline, to show comparative accuracy versus a dataset with no style information. Figure 8 displays captions generated using this method.

6.3 Text based style retrieval (text2style)

We explore the potential of the ALADIN-ViT+CLIP model trained in subsec. 6.1 to perform image retrieval, using textual tag queries. By first indexing the score assigned to each tag in the dictionary at the image level, we can then use a tag query to retrieve images based on the sorted scores for that tag. We use nearest-neighbour search using the image embeddings, reversing the tags generation experiment. Fig 9 shows some example image retrievals using text queries.

To measure the quality of the results, we run all text tags as queries. For each, we compute the WordNet similarity of the query text tag to the *k*th top tag associated with the image, following a tag retrieval using a given image. We vary *k* and collect the average scores at values of 1, 5, 10, and 25. The scores at these values are 0.72, 0.467, 0.392, and 0.332, respectively.

7 Conclusion

603

604

605

616

617

618

We proposed StyleBabel, a novel unique dataset of digital artworks and associated text describing their fine-grained artistic style. Our annotation approach was inspired by Grounded Theory (GT) [44], to support the 'emergence' of themes from the corpus of digital artwork – as opposed to fitting images to pre-existing categories [6]. These sessions generated discussion while simultaneously evolving and arriving through consensus at a shared vocabulary for describing image clusters of similar style.

We extended ALADIN [41] to incorporate a visual transformer [12] (ALADIN-ViT) encoder, obtaining state of the art style similarity discrimination, leveraging StyleBabel for the automated description of artwork images using keyword tags and captions. We also showed text-based image retrieval of images based on the generated style tags. Further work could explore use of tags as priors in generating captions, and exploring more downstream tasks using StyleBabel.

References	630
	631
1. Luis A. and Laura D. Esp: Labeling images with a computer game. pages 91–98, 01 2005.	7 632
2. P. Achiloptas, M. Ovsjanikov, K. Haydarov, M. Elnoseiny, and L. Guibas. Artemis: Affectiv	e 633
3 Zechen Bai Yuta Nakashima and Noa Garcia Explain me the painting Multi-tonic know	634
edgeable art description generation. <i>CoRR</i> , abs/2109.05743, 2021. 1	635
4. S. Bell and K. Bala. Learning visual similarity for product design with convolutional neura	al <u>636</u>
networks. In Proc. ACM SIGGRAPH, 2015. 3	637
5. Simone Bianco, Luigi Celona, Paolo Napoletano, and Raimondo Schettini. Predicting imag	e 620
aesthetics with deep learning. volume 10016, pages 117-125, 10 2016. 3	000
6. Kathy C. Constructing Grounded Theory : A Practical Guide through Qualitative Analysi	s. 639
Sage, London, 2006. 2, 4, 5, 14	640
<i>In Proc. FLMAP</i> 2012 2	n. 641
111 Proc. LLMAR, 2015. 5 8 T. Chen S. Kornblith M. Norouzi and G. Hinton A simple framework for contrastiv	642
learning of visual representations arXiv preprint arXiv:2002.05709.2020.10	643
9. J. Collomosse, T. Bui, M. Wilber, C. Fang, and H. Jin. Sketching with style: Visual search	h 644
with sketches and aesthetic context. In Proc. ICCV, 2017. 3	645
10. M. Cornia, M. Stefanini, L. Baraldi, and R. Cucchiara. Meshed-memory transformer for	or 646
image captioning. arXiv preprint arXiv:1912.08226, 2020. 1, 4	647
11. Karan Desai and Justin Johnson. Virtex: Learning visual representations from textual annota	a- 047
tions. CoRR, abs/2006.06666, 2020. 4, 14	648
12. A. DOSOVIISKIY, L. DEVEL, A. KOLESHIKOV, D. WEISSENDOHI, A. Zhai, I. Unterthinder, M. Do	3- 649 h
16x16 words: Transformers for image recognition at scale arXiv preprint arXiv:2010.1102	0 650
2020 9 10 14	, 651
13. Ali F., Mohsen H., M. Amin S., Peter Y., Cyrus R., Julia H., and David F. Every pictur	e 652
tells a story: Generating sentences from images. In Kostas D., Petros M., and Nikos F	653
editors, Computer Vision - ECCV 2010, pages 15-29, Berlin, Heidelberg, 2010. Springe	er 654
Berlin Heidelberg. 4	655
14. Eric F. Art History and its Methods: A critical anthology. Phaidon, London, 1995. 1, 2, 5	656
15. C. Fellbaum. WordNet: An Electronic Lexical Database. Bradford Books, 1998. 11	b 657
multi modal ratriaval CoPP abs/1810.00617, 2018, 1, 3	n 057
17 I. A Gatys A S Ecker and M Bether A neural algorithm of artistic style <i>arXiv preprin</i>	11
arXiv: 1508.06576, 2015. 2. 3	659
18. G. Ghiasi, H. Lee, M. Kudlur, V. Dumoulin, and J. Shlens. Exploring the structure of a rea	.l- ⁶⁶⁰
time, arbitrary neural artistic stylization network. arXiv preprint arXiv:1705.06830, 201	7. <mark>66</mark> 1
3	662
19. Koustav Ghosal, Aakanksha Rana, and Aljosa Smolic. Aesthetic image captioning from	n 663
weakly-labelled photographs. <i>CoRR</i> , abs/1908.11310, 2019. 3	664
20. Lun Huang, wenmin wang, Jie Chen, and Xiaoyong wei. Attention on attention for imag	e 665
21 X Huang and S Belongie Arbitrary style transfer in real-time with adaptive instance no	r- 666
malization In Proc ICCV 2017 3	I- 0000
22. X. Huang, MY. Liu, S. Belongie, and J. Kautz. Multimodal unsupervised image-to-imag	e
translation. In ECCV, 2018. 3	668
23. J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and supe	r- ⁶⁶⁹
resolution. In Proc. ECCV, 2016. 3	670
24. S. Karayev, M. Trentacoste, H. Han, A. Agarwala, T. Darrell, A. Hertzmann, and H. Win	n- 671
nemoller. Recognizing image style. In <i>Proc. BMVC</i> , 2014. 3	672
25. U. Kulkarili, V. Premiraj, V. Ordonez, S. Dhar, S. Li, Y. Choi, A. C. Berg, and I. L. Berg, Bobytolk: Understanding and generating simple image descriptions. <i>IEEE Transactions of the second sec</i>	g. 673
Pattern Analysis and Machine Intelligence 35(12):2801_2003 2013 4	674
1 anom 11mm you and machine memory 55(12).2091–2905, 2015, T	

 675
 26. Alon Lavie and Abhaya Agarwal. Meteor: An automatic metric for mt evaluation with high levels of correlation with human judgments. pages 228–231, 07 2007. 14
 676

 676
 27. Yu Lei, Albert Marrão Drauda, Thicker a Human and F. Y. Harradan, An anticher and Human and F. S. Human and Human and

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

- 27. Xu Lei, Albert Meroño-Peñuela, Zhisheng Huang, and F. V. Harmelen. An ontology model for narrative image annotation in the field of cultural heritage. In *WHiSe@ISWC*, 2017. 1
 28. X. Li, X. Yin, C. Li, D. Zhang, Y. Liu, J. Zhang, H. Hu, L. Dung, F. Wei, Y. Chailan, J. Chailan, J. S. Katalan, J. S.
- 28. X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei, Y. Choi, and J. Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. *arXiv* preprint arXiv:2004.06165, 2020. 1, 4
- 29. Y. Li, C. Fang, J. Yang, Z. Wang, X. Lu, and M. Yang. Universal style transfer via feature transforms. In *Proc. NIPS*, 2017. 3
- 30. Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. page 10, 01 2004.
 14
- 31. Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick,
 James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft
 COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. 14
- 32. K. Pang, Y. Yang, T. M. Hospedales, T. Xiang, and Y. Song. Solving mixed-modal jigsaw
 puzzle for fine-grained sketch-based image retrieval. In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10344–10352, Los Alamitos, CA, USA, jun 2020. IEEE Computer Society. 1
- 33. Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, page 311–318, USA, 2002. Association for Computational Linguistics. 14
- 34. T. Park, J.-Y. Zhu, O. Wang, J. Lu, E. Shechtman, A. Efros, and R. Zhang. Swapping autoen coder for deep image manipulation. In *Proc. ECCV*, 2020. 3
- 35. Andrea Pinotti. Formalism and the History of Style, pages 75 90. Brill, Leiden, The Netherlands, 2012. 2
- A. Radford, J. Wook Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell,
 P. Mishkin, J. Clark, G. Krueger, and I. Sutskever. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 3, 4, 10, 11, 12, 13
- 69937. A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever.700Zero-shot text-to-image generatio. arXiv preprint arXiv:2102.12092, 2021. 4
- 38. A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever.
 Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*, 2021. 4
- 39. Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time
 object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015. 14
- 40. M. Ruder, A. Dosovitskiy, and T. Brox. Artistic style transfer for videos. In *Proc. GCPR*, 2016. 3
- 41. D. Ruta, S. Motiian, B. Faieta, Z. Lin, H. Jin, A. Filipkowski, A. Gilbert, and J. Collomosse.
 Aladin: All layer adaptive instance normalization for fine-grained style similarity. *arXiv* preprint arXiv:2103.09776, 2021. 2, 3, 6, 9, 10, 14
- 42. Babak Saleh and Ahmed Elgammal. Large-scale classification of fine-art paintings: Learning the right metric on the right feature, 2015. 1
- 43. L. Shamir, T. Macura, N. Orlov, and D. Eckley. Impressionism, expressionism, surrealism: automated recognition of painters and schools of art. *IEEE Trans. Applied Perception*, 2010.
 3
- 44. J. Simondsen and T. Roberton. *Routledge International Handbook of Participatory Design*.
 Routledge, 2013. 2, 4, 5, 14
- 45. K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. **2**, **3**
- 46. A. Srinivas, T.-Yi Lin, N. Parmar, J. Shlens, P. Abbeel, and A. Vaswani. Bottleneck trans formers for visual recognition. *arXiv preprint arXiv:2101.11605*, 2021. 2
- 47. D. Ulyanov, V. Lebedev, A. Vedaldi, and V. Lempitsky. Texture networks: Feed-forward synthesis of textures and stylized images. In *Proc. ICML*, 2016. 3 3718 Synthesis of textures and stylized images. In *Proc. ICML*, 2016. 3
- 48. Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based
 image description evaluation. *CoRR*, abs/1411.5726, 2014. 14

720	49.	O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption	720
721	50	generator. arXiv preprint arXiv:1411.4555, 2015. 4	721
722	50.	X. Wang, G. Oxholm, D. Zhang, and YF. Wang. Multimodal transfer: A hierarchical deep	722
723	51.	XS. Wei, JH. Luo, J. Wu, and ZH. Zhou. Selective convolutional descriptor aggregation	723
724		for fine-grained image retrieval. arXiv preprint arXiv:1604.04994, 2017. 1	724
725	52.	M. J. Wilber, C. Fang, H. Jin, A. Hertzmann, J. Collomosse, and S. Belongie. Bam!	725
726		the behance artistic media dataset for recognition beyond photography. <i>arXiv preprint</i>	726
727	53.	Lei Xu and Xiaoguang Wang. Semantic description of cultural digital images: Using a	727
728	00.	hierarchical model and controlled vocabulary. <i>D Lib Mag.</i> , 21(5/6), 2015. 1	728
729	54.	B. Yao, A. Khosla, and L. Fei-Fei. Combining randomization and discrimination for fine-	729
730	55	grained image categorization. In CVPR 2011, pages 1577–1584, 2011. 1 D. Zhang, X. Li, X. Hu, J. Yang, L. Zhang, L. Wang, Y. Choi, and J. Cao, Vinul, Paviaiting	730
731	55.	visual representations in vision-language models. <i>arXiv preprint arXiv:2101.00529.2021.</i>	731
732		4	732
733	56.	L. Zhou, H. Palangi, L. Zhang, H. Hu, J. J. Corso, and J. Gao. Unified vision-language	733
734	57	pre-training for image captioning and vqa. arXiv preprint arXiv:1909.11059, 2019. 4	734
735	57.	genre: an analysis of features and classifiers. In <i>Proc. IEEE Workshop on Multimedia Sig</i>	735
736		Proc. (MMSP), 2009. 3	736
737			737
738			738
739			739
740			740
741			741
742			742
743			743
744			744
745			745
746			746
747			747
748			748
749			749
750			750
751			751
752			752
753			753
754			/54
755			755
756			756
757			757
750			750
760			759
761			761
762			760
763			762
764			764
704			104