GENERALIZATION AND KNOWLEDGE TRANSFER IN ABSTRACT VISUAL REASONING MODELS

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

Paper under double-blind review

ABSTRACT

We study generalization and knowledge reuse capabilities of deep neural networks in the domain of abstract visual reasoning (AVR), employing Raven's Progressive Matrices (RPMs), a recognized benchmark task for assessing AVR abilities. Two knowledge transfer scenarios referring to the I-RAVEN dataset are investigated. Firstly, inspired by generalization assessment capabilities of the PGM dataset and popularity of I-RAVEN, we introduce Attributeless-I-RAVEN, a benchmark with 10 generalization regimes that allow to test generalization of abstract rules applied to held-out attributes. Secondly, we construct *I-RAVEN-Mesh*, a dataset that enriches RPMs with a novel component structure comprising line-based patterns, facilitating assessment of progressive knowledge acquisition in transfer learning setting. The developed benchmarks reveal shortcomings of the contemporary deep learning models, which we partly address with *Pathways of Normalized Group Convolution (PoNG)* model, a novel neural architecture for solving AVR tasks. PoNG excels in both presented challenges, as well as the standard I-RAVEN and PGM setups. Encouraged by these promising results, we further evaluate PoNG in another AVR task, visual analogy problem with both synthetic and real-world images, demonstrating its strength beyond PRMs.

028 1 INTRODUCTION

Generalization, the ability of a model to perform well on unseen 031 data, remains a fundamental challenge in deep learning (DL). While DL methods have demonstrated remarkable achievements in vari-033 ous domains, their generalization capabilities are often questioned, 034 particularly in tasks that demand abstract problem-solving and reasoning skills (Chollet, 2019). One such domain is abstract visual reasoning (AVR) (Mitchell, 2021; van der Maas et al., 2021; Stabinger et al., 2021; Małkiński & Mańdziuk, 2023) that encom-037 passes tasks requiring (human) fluid intelligence – an aspect of human cognition believed to be crucial for reasoning in neverencountered settings (Snow et al., 1984; Carpenter et al., 1990). 040 The most popular AVR tasks are Raven's Progressive Matrices 041 (RPMs) (Raven, 1936; Raven & Court, 1998), which constitute a 042 common problem found in human IQ tests. Typical RPMs com-043 prise two components – the context panels arranged in a 3×3 grid 044 with the bottom-right panel missing and up to 8 answer panels, out



Figure 1: **RPM example.** The correct answer is A.

of which only one correctly completes the matrix. Solving an RPM instance requires identification
 of underlying abstract rules applied to certain attributes of the objects composing the instance (see
 Fig. 1 for an illustrative example).

Design of computational methods capable of tackling RPMs has for decades been an active area of research (Evans, 1964; Gentner, 1980; Foundalis, 2006; Lovett et al., 2007; Kunda et al., 2010; Strannegård et al., 2013). Consequently, a number of works considered automatic creation of RPM datasets (Matzen et al., 2010; Wang & Su, 2015; Mańdziuk & Żychowski, 2019) and a wide suite of predictive models (Hernández-Orallo et al., 2016; Hernández-Orallo, 2017) were proposed, with DL methods showing the most promising performance (Yang et al., 2022; Małkiński & Mańdziuk, 2022). While this rapid progress led to exceeding the human level in particular problem setups (Wu

et al., 2020; Mondal et al., 2023), a fundamental challenge of generalization to novel problem settings remains largely unattained.

Initial works designed several RPM datasets (Matzen et al., 2010; Wang & Su, 2015; Hoshen & 057 Werman, 2017), however, measuring generalization was not their focus. While some works explored knowledge transfer between related tasks (Mańdziuk & Zychowski, 2019; Tomaszewska et al., 2022), the complexity of the datasets was limited and consequently they didn't pose a chal-060 lenge for contemporary DL methods. To measure generalization in modern DL models, the PGM 061 dataset was introduced (Barrett et al., 2018). PGM defines eight generalization regimes, each spec-062 ifying the distribution of objects, rules and attributes in train and test splits. For instance, in the 063 Held-out Triples split, a given rule-object-attribute triplet (e.g. Progression on Object's 064 Size) was assigned only to one of the two splits. In effect, the models were tested on triplet combinations different from training ones, allowing to assess their generalization capabilities. A subsequent 065 work proposed RAVEN (Zhang et al., 2019a), another RPM dataset with enriched perceptual com-066 plexity of matrices instantiated in seven visual configurations (Center, 2x2Grid, 3x3Grid, 067 Left-Right, Up-Down, Out-InCenter, Out-InGrid). Moreover, the benchmark is char-068 acterized by a moderate sample size, i.e. 70K instances, compared to 1.42M RPMs per each of the 069 eight regimes in PGM. Due to this size disparity, subsequent research gravitated towards RAVEN and its revised variants (I-RAVEN (Hu et al., 2021) and RAVEN-Fair (Benny et al., 2021)), which 071 didn't require substantial computational resources to train DL models. 072

Drawing inspiration from the broad adoption of RAVEN and the generalization assessment capabilities of PGM, this paper proposes a novel suite of generalization challenges stemming from I-RAVEN (Hu et al., 2021) (a revised variant of RAVEN that removes a bias in RAVEN's answer panels). However, unlike I-RAVEN, the proposed suite of benchmarks allows for a direct assessment of the generalization and knowledge transfer of AVR models. Compared to PGM, our datasets feature compositionality and variety of figure configurations, and their processing doesn't require substantial computational resources. Furthermore, they include structural annotations, which are utilized, for example, in recent neuro-symbolic approaches.

First, we introduce *Attributeless-I-RAVEN*, comprising 10 generalization regimes. The 4 primary regimes correspond to specific held-out attributes ({Position, Type, Size, Color}), resp. The training matrices in these regimes adhere to the Constant rule for the respective attribute, whereas test matrices employ a rule different from Constant for this attribute (i.e., Progression, Arithmetic, or Distribute Three). Moreover, we propose 6 extended regimes: 3 of them feature a held-out attribute pair, while another 3 replace the Constant rule in the training set with each remaining rule. In effect, each regime comprises different distributions of training and test data.

Next, we propose *I-RAVEN-Mesh*, a variant of I-RAVEN with a new grid-like structure overlaid on the matrices. The dataset enables assessing generalization to incrementally added structures and progressive knowledge acquisition in a transfer learning (TL) setting.

In investigations involving 13 contemporary AVR DL models, we observed that the introduced benchmarks present a substantial challenge for the tested methods. This prompted the development of *Pathways of Normalized Group Convolution (PoNG)*, a novel AVR model that excels in both problem setups: generalization to held-out attributes and incremental knowledge acquisition.

- 096 Our main contributions can be summarized as follows:
- 097 098

099

102

- We introduce the *Attributeless-I-RAVEN* (A-I-RAVEN) dataset that enables measuring generalization across 10 regimes.
- We construct *I-RAVEN-Mesh*, an extension of I-RAVEN with a new component structure that facilitates assessment of progressive knowledge acquisition in a TL setting.
- We evaluate the performance of state-of-the-art AVR models on the introduced benchmarks, uncovering their limitations in terms of generalization to novel problem settings.
- We propose a new neural architecture for solving AVR tasks termed PoNG, which excels in addressing both introduced challenges, as well as the standard I-RAVEN and PGM setups. Additionally, PoNG demonstrates the state-of-the-art performance in visual analogy problem (VAP) in both synthetic and real-world setups.

108 2 RELATED WORK

110 Generalization in AVR. In recent years, a variety of AVR problems and corresponding datasets 111 have emerged (Bongard, 1968; Nie et al., 2020; Fleuret et al., 2011; Qi et al., 2021; Shanahan et al., 112 2020; Jiang et al., 2024; Hill et al., 2019; Zhang et al., 2020) and several attempts have been made 113 to measure generalization in contemporary AVR models based on the introduced benchmarks. In particular, distinct visual configurations were employed in RAVEN to assess how a model trained on 114 one configuration performs on the remaining ones (Zhang et al., 2019a; Spratley et al., 2020; Zhuo 115 & Kankanhalli, 2021). Although in such a setting the visual aspects of train/test matrices come from 116 different distributions, the underlying rules and attributes remain the same. In contrast, A-I-RAVEN 117 enables studying the generalization of rules applied to held-out attributes, shifting the focus from 118 perception towards reasoning. Besides RPMs, the limits of generalization have been explored in 119 other AVR tasks as well. Visual Analogy Extrapolation Challenge evaluates model's capacity for 120 extrapolation (Webb et al., 2020). However, such specialized datasets might favor models that ex-121 plicitly embed the notion of extrapolation in their design and aim for being invariant only to specific 122 attributes such as object size or location. Differently, our benchmarks allow verifying the model's 123 capacity to learn a given concept from the data and generalize it to novel settings. This perspective 124 links our work to the recent literature on concept learning (Odouard & Mitchell, 2022; Moskvichev 125 et al., 2023). However, the concept-oriented benchmarks that originate from ARC (Chollet, 2019) remain largely unsolved by DL models and pose a significant challenge even for leading multi-modal 126 large language models (Mitchell et al., 2023). In contrast, both benchmarks proposed in this work 127 are attainable by DL models, though further advances in generalization abilities of the models are 128 necessary to consider them solved. 129

130 Model architectures. Preliminary attempts to solve RPMs with DL models involve WReN (Bar-131 rett et al., 2018) that reasons over object relations using Relation Network (Santoro et al., 2017), 132 or SRAN (Hu et al., 2021) that relies on a hierarchical architecture with panel encoders devoted 133 to particular image groups. A common theme enabling generalization in DL models is to explicitly 134 identify RPM objects. To this end, RelBase (Spratley et al., 2020) employs Attend-Infer-Repeat (Es-135 lami et al., 2016), an unsupervised scene decomposition method, STSN (Mondal et al., 2023) utilizes 136 Slot attention (Locatello et al., 2020) to decompose matrix to slots containing particular objects and 137 Temporal Context Normalization (TCN) (Webb et al., 2020) to normalize latent matrix panel representations in a task-specific context, DRNet (Zhao et al., 2024) relies on a dual-stream design, and 138 MRNet (Benny et al., 2021) presents a multi-scale architecture. SCL (Wu et al., 2020) proposes 139 the scattering transformation, CoPINet (Zhang et al., 2019b) and CPCNet (Yang et al., 2023b) rely 140 on contrastive architectures, PredRNet (Yang et al., 2023a) learns to minimize the prediction error, 141 ALANS (Zhang et al., 2021) and PrAE (Zhang et al., 2022a) employ neuro-symbolic architectures, 142 and SCAR (Małkiński & Mańdziuk, 2024b) adapts its computation to the structure of the consid-143 ered matrix. Despite the high variety of AVR models, experiments on the introduced benchmarks 144 reveal their shortcomings in terms of generalization and knowledge transfer. In this context, we pro-145 pose PoNG, a new AVR model that excels in the presented tasks by combining parallel architecture, 146 weight sharing, and tactical normalization.

3 Methods

147 148

149

155 156

The set of attributes in I-RAVEN is $\mathcal{A} = \{\text{Position, Number, Type, Size, Color}\}$ and the set of rules is $\mathcal{R} = \{\text{Constant, Progression, Arithmetic, Distribute Three}\}$. For attribute $a \in \mathcal{A}$ and a dataset split $s \in \mathcal{S}$, where $\mathcal{S} = \{\text{train, val., test}\}$, we define the set of rules applicable to a in split s by $R(a, s) \subseteq \mathcal{R}$. In I-RAVEN all rule-attribute pairs are valid in all splits:

$$R(a,s) = \mathcal{R}, \quad \forall a \in \mathcal{A} \land \forall s \in \mathcal{S}$$
(1)

157 3.1 ATTRIBUTELESS-I-RAVEN

To probe generalization in DL models, we present A-I-RAVEN, a benchmark composed of 10 generalization regimes. Example matrices are illustrated in Fig. 2, with additional samples provided in Appendix A. Each regime defines a set of held-out attributes A^* , each with a corresponding rule $r^*(a), a \in A^*$. In train and validation splits, held-out attribute $a \in A^*$ is governed by $r^*(a)$. In



Figure 2: Attributeless-I-RAVEN. Left: Matrices from the A/Position regime belonging to the 2×2 Grid configuration. In (a), object position is constant across rows, while in (b) object numerosity is governed by Distribute Three. Right: Matrices from the A/Color regime belonging to the Left-Right configuration. In (c), object color is constant across rows in left and right image parts, while in (d) it's governed by Progression. Correct answers are marked in a green dotted border. Please refer to Appendix A for examples from other generalization regimes.

the test split, $a \in A^*$ is governed by a different rule sampled from $\mathcal{R} - \{r^*(a)\}$. In effect, during training, the model doesn't see rule–attribute combinations required to solve test matrices. There are no rule-related constraints on the remaining attributes. In summary, we have:

$$R(a,s) = \begin{cases} \{r^*(a)\} & \text{if } a \in A^* \land s \in \{\text{train, validation}\}, \\ \mathcal{R} - \{r^*(a)\} & \text{if } a \in A^* \land s = \text{test}, \\ \mathcal{R} & \text{if } a \notin A^*. \end{cases}$$
(2)

We define 4 primary regimes with $r^*(a) = \text{Constant}$ that correspond to individual held-out at-189 tributes ($|A^*| = 1$), denoted as A/<Attribute> (e.g., A/Type). Since Position and Number 190 attributes are tightly coupled (e.g., it's impossible to increase cardinality of objects while keeping 191 their position constant), we allocate a single generalization regime, A/Position, to cover both 192 attributes. In addition, we define 6 extended regimes as supplementary generalization challenges. 193 In the first group a pair of attributes is held-out in the training set, i.e. $|A^*| = 2$. Specifically, 194 we introduce 3 new regimes: A/ColorSize, A/ColorType, and A/SizeType, based on the 195 respective attribute pairs. In the second group, Constant rule in $r^*(a)$ is replaced with each 196 of the 3 remaining rules, leading to A/Color-Progression, A/Color-Arithmetic, and 197 A/Color-DistributeThree regimes. While this modification could be applied to all the de-198 scribed regimes, we focus on the Color attribute due to its broad range of possible values.

200 3.2 I-RAVEN-MESH

180 181

182

183

199

201

The other of the proposed benchmarks is designed to probe progressive knowledge acquisition in a 202 TL setting. I-RAVEN-Mesh extends I-RAVEN by introducing a novel visual component overlaid 203 on top of the existing I-RAVEN components (see Fig. 3). Though the dataset can serve as a learning 204 challenge on its own, the main motivation behind its introduction is to employ models pre-trained on 205 I-RAVEN and fine-tune them on I-RAVEN-Mesh with a configurable train sample size, facilitating 206 analysis of their TL performance. The mesh grid comprises from 1 to 12 lines placed in predefined 207 locations. The set of available lines covers the inner and outer edges of a 2×2 grid (12 lines in 208 total). The mesh component has two attributes: $\mathcal{A}^{\text{mesh}} = \{\text{Number}, \text{Position}\}, \text{ which govern}\}$ 209 the count and location of lines, respectively. To each attribute a rule $r \in \mathcal{R}$ can be applied. Table 1 210 describes the effect of applying a given rule-attribute pair to the mesh component. To generate 211 the mesh component of an I-RAVEN-Mesh matrix, we sample one of the two attributes $a \in \mathcal{A}^{\text{mesh}}$ 212 and a corresponding rule $r \in \mathcal{R}$ that governs its values. As the attributes often depend on each other 213 (e.g., it's impossible to increase the number of lines while keeping their position constant), we don't constrain the value of the other attribute. The rule-attribute pairs for the base I-RAVEN components 214 are generated in the same way as in the original dataset. To generate answers to the matrix, we follow 215 the impartial algorithm proposed in I-RAVEN (Hu et al., 2021). In addition, each matrix contains



Figure 3: I-RAVEN-Mesh. Matrices with the Position attribute of the mesh component governed by all applicable rules. For the sake of readability, we present examples belonging to the Center configuration. (a) Line position is constant in each row. (b) The line pattern displayed in the first column is rotated by 90 degrees in subsequent columns. (c) The union set operator applied to the first and the second column produces line positions in the third column. (d) Each row contains lines arranged in one out of three available patterns. Correct answers are marked in a green dotted border. Please refer to Appendix A for examples concerning the Number attribute.

Table 1: Description of rule-attribute pairs in I-RAVEN-Mesh. D3 marks Distribute Three.

Attribute	Rule	Description
	Constant	Each image in a given row contains the same number of lines.
	Progression	The count of lines in a given row changes by a constant factor (e.g. 2, 4, 6).
Number	Arithmetic	The number of lines in the third column is determined based on an arithmetic
	_	operation applied to the preceding columns (e.g. $3 - 1 = 2$).
	D3	Three line counts are sampled and spread among images in a given row.
	Constant	Each image in a given row contains the same position of lines.
	Progression	A panel arrangement is sampled in each row and rotated by 90 degrees in
Position		subsequent columns.
POSICION	Arithmetic	The position of lines in the third column is computed based on a set operation
		(union or difference) applied to the preceding columns.
	D3	Three line arrangements are sampled and spread among images in a given row.

at least one incorrect answer that differs from the correct one only in the mesh component, ensuring that the solver has to identify the correct rule governing the mesh component in order to solve the matrix. To facilitate training with an auxiliary loss, in which the model additionally predicts the representation of rules governing the matrix (Barrett et al., 2018), we extend the base set of rule annotations with ones concerning the Mesh component.

3.3 PATHWAYS OF NORMALIZED GROUP CONVOLUTION (PONG)

259 In initial experiments, we've found out that SOTA AVR models struggle in the proposed generaliza-260 tion challenges. Consequently, we introduce PoNG (Fig. 4), a novel model that outcompetes baselines across a number of problem settings. The model follows a typical two-stage design. Firstly, it 261 generates an embedding of each image panel. Then, it aggregates representations of matrix panels 262 to predict the index of the correct answer. The details are described in Appendix D. 263

264 Let (X, y, r) denote an RPM, where $X = \{x_i\}_{i=1}^{16}$ is the set of image panels comprising 8 context panels $\{x_i\}_{i=1}^8$ and 8 answer panels $\{x_i\}_{i=9}^{16}$, $x_i \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$ is the one-hot encoded index of the correct answer, $r \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$ is the one-hot encoded index of the correct answer, $r \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$ is the one-hot encoded index of the correct answer, $r \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$ is the one-hot encoded index of the correct answer, $r \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$ is the one-hot encoded index of the correct answer, $r \in [0, 1]^{h \times w}$, i = 1, ..., 16 is a grayscale image of height h and width $w, y \in \{0, 1\}^8$. 265 266 $\{0,1\}^{d_r}$ is the multi-hot encoded representation of matrix rules of dimensionality d_r using sparse 267 encoding (Małkiński & Mańdziuk, 2024a). In each experiment h = w = 80, while d_r is determined 268 by the number of matrix components in the corresponding dataset ($d_r = 48$ for I-RAVEN-Mesh, 269 $d_r = 40$ otherwise; see Appendix C for details).

250

253

254

255 256 257

258

228

229

230

231

232

233

234 235

236



288 Figure 4: **PoNG.** (a) The panel encoder embeds each input image x_i independently, producing h_i . 289 Context panel embeddings ${h_i}_{i=1}^8$ together with the embedding of k'th answer h_k are stacked and processed with the reasoner, leading to z_k . (b) The pathways block, a key component of PoNG, 290 comprises four parallel pathways P1 – P4. (c) P3 and (d) P4 employ novel normalized group convolution operators. PosEmb denotes position embedding, G-C the group convolution module used in 292 P3, and GP-C the group-pair convolution module used in P4. The red dashed line marks the point 293 after which G-C and GP-C perform analogous computation.

296 **Panel encoder.** The first component of the model has the form $\mathcal{E}: x \to h$, where $h \in \mathbb{R}^{d_h}$ is the 297 input panel embedding of dimensionality d_h . Following RelBase (Spratley et al., 2020), the module 298 comprises 2 blocks of the same architecture. Each block includes 2 parallel pathways that build 299 high-level and low-level features, resp. The first one contains 2 convolutional blocks, each with 300 2D convolution, ReLU, and Batch Normalization (BN) (Ioffe & Szegedy, 2015). The second one 301 contains 2D max pooling followed by 2D convolution. The sum of both pathway results forms the 302 block output. Differently from RelBase, we flatten the height and width dimensions of the resultant embedding, pass it through a linear layer with ReLU, flatten the channel and spatial dimensions, 303 and pass the tensor through a feed-forward residual block with Layer Normalization (LN) (Ba et al., 304 2016). Finally, we concatenate the tensor with a position embedding (a learned 25-dimensional 305 vector for each cell in the 3×3 context grid), leading to h. 306

Reasoner. The second component of the model has the form $\mathcal{R}: \{h_i\}_{i=1}^8 \cup h_k \to z_k$, where h_k 308 is the panel embedding of k'th answer. For each answer panel, the reasoner produces embedding 309 z_k that describes how well the considered answer fits into the matrix context. Panel embeddings 310 $\{h_i\}_{i=1}^{8} \cup h_k$ are stacked and processed by a sequence of 3 reasoning blocks interleaved with 2 311 bottleneck layers for dimensionality reduction. Each reasoning block comprises BN and 4 parallel 312 pathways, outputs of which are added together to form the output of the block. Next, the latent 313 representation is passed through adaptive average pooling, flattened, processed with a linear layer 314 with ReLU, passed through BN and projected with a linear layer to $z_k \in \mathbb{R}^{128}$.

315 316

307

291

295

Pathways. The key aspect of the reasoner module are its pathways. Each takes an input tensor 317 of shape (B, C, D), where B is the batch size, C is the number of channels, and D is the feature 318 dimension. In the first reasoning block $D = d_h$ and C = 9 corresponds to the number of panel 319 embeddings in the considered group. Pathways are described as follows: P1 - a pointwise 1D 320 convolution layer that mixes panel features at each spatial location; P2 - a sequence of 2 blocks, each 321 comprising 1D convolution, ReLU, and BN, that builds higher level features spanning neighbouring spatial locations; P3 – analogous to P2, but 1D convolution is replaced with a group 1D convolution 322 that splits the tensor into several groups along the channel dimension, applies a 1D convolution with 323 shared weights to each group, and adds together the representations of each group; P4 – analogous

to P3, but groups are arranged into pairs concatenated along the channel dimension and processed with a 1D convolution with shared weights. In contrast to (Krizhevsky et al., 2012), the proposed group convolution layers in both P3 and P4 apply TCN (Webb et al., 2020) to the outputs in each group. In the first layer P3 and P4 split the input tensor into 3 groups, which allows for producing embeddings of each matrix row and each pair of rows, resp. Though we apply the pathways block in the RPM context, we envisage it as a generic module, also applicable to other settings involving a set of vector representations of shape (B, C, D).

Answer prediction. Eight representations of the context matrix filled-in with the respective an-332 swer, $\{z_k\}_{k=1}^8$, are processed with three prediction heads. The target head $\mathcal{P}^y: z_k \to \widehat{y_k}$ employs 333 two linear layers interleaved with ReLU to produce score $\widehat{y_k} \in \mathbb{R}$ describing how well the an-334 swer k aligns with the matrix context. The aggregate rule head $\mathcal{P}_1^r : \{z_k\}_{k=1}^8 \to \hat{r_1}$ computes 335 the sum of inputs and processes it with two linear layers interleaved with ReLU, producing a la-336 tent prediction of matrix rules $\hat{r}_1 \in \mathbb{R}^{d_r}$. The target-conditioned rule head $\mathcal{P}_2^r : \{z_k\}_{k=1}^8 \to \hat{r}_2$ 337 processes its input through a linear layer and computes a weighted sum of the resultant embed-338 dings with weights given by the predicted probability distribution over the set of possible an-339 swers $\sigma(\{\widehat{y}_k\}_{k=1}^8)$, where σ denotes softmax. The model is trained with a joint loss function $\mathcal{L} = \operatorname{CE}(\sigma(\{\hat{y}_k\}_{k=1}^{3}), y) + \beta \operatorname{BCE}(\zeta(\hat{r_1}, r)) + \gamma \operatorname{BCE}(\zeta(\hat{r_2}, r)), \text{ where } \zeta \text{ denotes sigmoid, CE cross-entropy, BCE binary cross-entropy, } \beta = 25 \text{ and } \gamma = 5 \text{ are balancing coefficients.}$ 340 341

342 343

344

349

331

4 EXPERIMENTS

We assess generalization of state-of-the-art models for solving RPMs on A-I-RAVEN, evaluate progressive knowledge acquisition on I-RAVEN-Mesh, and conduct an ablation study to showcase the contribution of the respective modules that constitute PoNG. We also evaluate PoNG on two additional VAPs comprising synthetic (Hill et al., 2019) and real-world (Bitton et al., 2023) images.

Experimental setup. In all experiments we use the Adam optimizer (Kingma & Ba, 2014) with 350 $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$ and a batch size set to 128. Learning rate is initialized to 0.001 and 351 reduced 10-fold (at most 3 times) if no progress is seen in the validation loss in 5 subsequent epochs, 352 and training stops early in the case of 10 epochs without progress. Unless stated otherwise, each 353 model configuration was trained 3 times with a different seed, and we report mean and standard 354 deviation for these runs. In each experiment, we utilize 42 000 training, 14 000 validation, and 355 14 000 test matrices, following the standard data split protocol taken in prior works (Zhang et al., 356 2019a; Hu et al., 2021). All reference models are trained with the auxiliary loss with sparse encoding 357 and $\beta = 1$. Experiments were run on a worker with a single NVIDIA DGX A100 GPU. 358

Baselines. In addition to the simple CNN-LSTM baseline (Barrett et al., 2018), we assess generalization of SOTA AVR models including WReN (Barrett et al., 2018), CoPINet (Zhang et al., 2019b),
RelBase (Spratley et al., 2020), SCL (Wu et al., 2020), MRNet (Benny et al., 2021), ALANS (Zhang
et al., 2021), SRAN (Hu et al., 2021), PrAE (Zhang et al., 2022a), CPCNet (Yang et al., 2023b), PredRNet (Yang et al., 2023a), STSN (Mondal et al., 2023), and DRNet (Zhao et al., 2024). For direct
comparison, we evaluate all models on I-RAVEN following the above-described experimental setup.

- Reproducibility. To guarantee reproducibility of experiments, we use a fixed set of random seeds
 and turn off hardware and framework features concerning indeterministic computation wherever
 possible. Together with the code, we provide the full training script that can be used to run all training jobs. The training job is packaged as a Docker image with fixed dependencies to isolate the configuration of the training environment. The released code allows for generation of all datasets from
 scratch, eliminating the dependency on file-hosting services required to distribute the data. The code
 for reproducing all experiments is publicly accessible at: <hidden-for-blind-review>.
- Generalization on Attributeless-I-RAVEN. In the first set of experiments we evaluate all considered models on 4 primary generalization regimes of A-I-RAVEN. The results are presented in Table 2, along with the reference results on I-RAVEN and I-RAVEN-Mesh. PoNG outperforms all selected baselines across all settings. Among baseline models, the best results on A/Color and A/Position are achieved by DRNet, followed by RelBase and SCL. In the remaining attributeless regimes, SCL outperforms other baselines with DRNet taking the second place. Interestingly,

I-KVIN ¹ denotes results on I-KAVEN reported by model authors in the corresponding papers.								
	I-RVN [†] $ $	I-RAVEN	Mesh	A/Color	A/Pos.	A/Size	A/Type	
ALANS	_	$27.0 (\pm 8.4)$	$15.9 (\pm 2.6)$	$15.2 (\pm 1.4)$	$16.0 (\pm 1.0)$	$23.3 (\pm 6.5)$	$19.0 (\pm 3.4)$	
CPCNet	98.5	$70.4 (\pm 6.4)$	$66.6 (\pm 5.1)$	$51.2 (\pm 3.8)$	$68.3 (\pm 4.0)$	$43.5 (\pm 3.5)$	$38.6 (\pm 4.3)$	
CNN-LSTM	18.9	$27.5 (\pm 1.5)$	$28.9(\pm 0.4)$	$17.0 (\pm 3.1)$	$24.0 (\pm 2.9)$	$13.6 (\pm 1.4)$	$14.5 (\pm 0.8)$	
CoPINet	46.1	$43.2 (\pm 0.1)$	$41.1 (\pm 0.3)$	$32.5 (\pm 0.2)$	$41.3 (\pm 1.6)$	$21.8 (\pm 0.2)$	$19.8 (\pm 0.9)$	
DRNet	97.6	$90.9(\pm 1.1)$	$83.9(\pm 2.7)$	$70.0 (\pm 1.6)$	$77.5 (\pm 0.9)$	$54.3 (\pm 3.0)$	$44.3 (\pm 0.8)$	
MRNet	83.5	$86.7 (\pm 2.3)$	$79.5 (\pm 2.0)$	$33.6 (\pm 8.2)$	$62.6 (\pm 2.6)$	$20.6 (\pm 5.0)$	$19.4 (\pm 0.3)$	
PrAE	77.0	$19.5 (\pm 0.4)$	$33.2 (\pm 0.4)$	$47.9 (\pm 0.9)$	$68.2 (\pm 3.3)$	$41.3 (\pm 1.8)$	$37.0 (\pm 1.7)$	
PredRNet	96.5	$88.8 (\pm 1.8)$	$59.2 (\pm 6.4)$	$59.4 (\pm 1.0)$	$73.7 (\pm 0.7)$	$47.5 (\pm 1.3)$	$40.2 (\pm 1.3)$	
RelBase	91.1	$89.6 (\pm 0.6)$	$84.9(\pm 4.4)$	$67.4 (\pm 2.7)$	$76.6 (\pm 0.3)$	$51.1 (\pm 2.4)$	$44.1 (\pm 1.0)$	
SCL	95.0	$83.4 (\pm 2.5)$	$80.9(\pm 1.5)$	$65.1 (\pm 2.0)$	$76.7 (\pm 7.1)$	$65.6 (\pm 2.4)$	$49.5(\pm 1.8)$	

 $57.8 (\pm 0.2)$

 $48.7 (\pm 11.5)$

 $25.7 (\pm 0.2)$

89.3 (± 2.4)

378 Table 2: Single-task learning. Mean and standard deviation of test accuracy for three random seeds. 379 Best dataset results are marked in bold and the second best are underlined. Pos. denotes Position.

Table 3: PoNG ablations. Test accuracy averaged across 3 random seeds and a difference to the
default model setup (cf. Table 2). Union denotes application of all ablations but the first one.

 $38.3 (\pm 1.0)$

 $39.3 (\pm 6.9)$

 $16.9 (\pm 0.5)$

80.3 (± 4.3)

 $56.9 (\pm 0.7)$

 $36.1 (\pm 19.9)$

 $17.3 (\pm 0.4)$

79.3 (± 0.7)

 $34.4 (\pm 3.0)$

 $38.4 (\pm 16.6)$

 $12.4 (\pm 0.5)$

73.5 (± 3.1)

 $30.7 (\pm 2.2)$

 $39.1 (\pm 5.0)$

 $15.1 (\pm 0.7)$

59.4 (± 6.9)

	I-RAVEN	Mesh	A/Color	A/Pos.	A/Size	A/Type
w/o P1 and P2 w/o P3 and P4	92.8(-3.1) 95.6(-0.3)	74.4(-14.9)	73.3(-7.0)	76.4(-2.9) 78.6(-0.7)	58.4(-15.2) 73.9(+.0.4)	49.5 (-9.8) 53 9 (-55)
w/o T S and T 4 w/o TCN	96.0 (-0.3) 96.0 (+0.1)	90.8 (+ 1.4)	75.4(-4.9)	80.3 (+ 1.0)	66.6(-6.9)	57.5(-1.9)
$egin{array}{ll} \gamma &= 0 \ eta &= 0 \end{array} \ eta &= 0 \end{array}$	95.7 (-0.1) 94.2 (-1.7)	$88.8 (-0.5) \\91.4 (+2.1)$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	79.6 (+ 0.3) 77.5 (- 1.8)	73.0(-0.5) 70.3(-3.2)	$56.9 (- 2.5) \\ 53.3 (- 6.1)$
$\begin{array}{l} \gamma = 0 \wedge \beta = 0 \\ \text{union} \end{array}$	$\begin{array}{c} 79.7 \ (-16.2) \\ 81.4 \ (-14.5) \end{array}$	$\begin{array}{c} 32.7 \ (-56.7) \\ 32.5 \ (-56.8) \end{array}$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{l} 75.1 \ (- \ 4.2) \\ 74.1 \ (- \ 5.2) \end{array}$	$\begin{array}{r} 64.9 \ (- \ 8.6) \\ 66.9 \ (- \ 6.6) \end{array}$	$\begin{array}{c} 49.0 \ (-10.3) \\ 46.0 \ (-13.4) \end{array}$

410 the top 3 models rely on rather shallow architectures, yet outcompete other methods that rely on a 411 deeper layout, such as SRAN or STSN. This suggests that parameter-efficient AVR models not only 412 excel in solving RPMs but also generalize better.

413 Generalization regimes of A-I-RAVEN pose a bigger 414 challenge than the base dataset. While PoNG, the best 415 performing model, achieved 95.9% test accuracy on I-416 RAVEN, on A-I-RAVEN regimes it scored from 59.4% 417 (on A/Type) to 80.3% (on A/Color). Fig. 5 displays the difference in PoNG's performance on test and 418 validation splits. On I-RAVEN and I-RAVEN-Mesh 419 the difference is negligible, as in these datasets both 420 splits follow the same distribution. However, the dif-421 ference in attributeless regimes is significant, which 422 indicates the need for further research on generaliza-423 tion. In Appendix E we present further evaluation on 424 6 extended A-I-RAVEN regimes. As shown in Ta-425 ble 10, replacing the Constant rule in the training set 426 with Progression or DistributeThree yields



Figure 5: Dataset difficulty. PoNG's performance on test and validation splits.

427 a dataset of similar complexity (the best model achieves 81.4 / 81.3% accuracy), while using the 428 Arithmetic rule increases the difficulty (the best model scored 70%). Furthermore, using a pair 429 of held-out attributes significantly increases the complexity. For instance, in A/SizeType, the most challenging regime, the best result is only 33.5%. Notably, PoNG outperforms all other mod-430 els in 5 out of 6 settings. We conclude that A-I-RAVEN provides a suite of challenging regimes of 431 variable complexity, in which even the best performing models are far from solving all test matrices.

8

392

393

394

396 397

SRAN

STSN

WReN

PoNG (ours)

60.8

95.7

23.8

95.9

 $58.2 (\pm 1.6)$

 $51.0 (\pm 24.8)$

 $18.4 (\pm 0.0)$

95.9 (± 0.7)

446

447

448 449

450

451

452

453



Figure 6: **Transfer learning.** Mean and standard deviation of test accuracy on I-RAVEN-Mesh across three random seeds. Models were trained in two setups: 1) from scratch on I-RAVEN-Mesh with variable sample size; 2) pre-trained on full I-RAVEN and fine-tuned on I-RAVEN-Mesh with variable sample size. Results for setups 1) and 2) are shown below and above the plot lines, resp.

Table 4: **PGM**. Test accuracy of PoNG in all regimes of the PGM dataset. The Interpolation regime is denoted as Inter., Held-out Attribute Pairs as HO-AP, Held-out Triple Pairs as HO-TP, Held-out Triples as HO-Triples, Held-out Attribute line-type as HO-LT, Held-out Attribute shape-colour as HO-SC, and Extrapolation as Extra. For reference, we provide results of SCL (Wu et al., 2020; Małkiński & Mańdziuk, 2024a), MRNet (Benny et al., 2021), ARII (Zhang et al., 2022b), PredR-Net (Yang et al., 2023a), DRNet Zhao et al. (2024), and Slot-Abstractor (Mondal et al., 2024).

Model	Neutral	Inter.	HO-AP	HO-TP	HO-Triples	HO-LT	HO-SC	Extra.	Avg.
SCL	87.1	56.0	79.6	76.6	23.0	14.1	12.6	19.8	46.1
MRNet	93.4	68.1	38.4	55.3	25.9	30.1	16.9	19.2	43.4
ARII	88.0	57.8	50.0	64.1	32.1	16.0	12.7	29.0	43.7
PredRNet	97.4	70.5	63.4	67.8	23.4	27.3	13.1	19.7	47.8
DRNet	99.1	83.8	93.7	78.1	48.8	27.9	13.1	22.2	58.3
Slot-Abstractor	91.5	91.6	63.3	78.3	20.4	16.7	14.3	39.3	51.9
PoNG (ours)	<u>98.1</u>	75.2	<u>92.1</u>	97.7	46.1	16.9	12.6	19.9	<u>57.3</u>

465 **Progressive knowledge acquisition on I-RAVEN-Mesh.** In the second set of experiments we 466 employ I-RAVEN-Mesh to examine the TL ability of the best performing models (see Appendix E 467 for extended results). To this end, we consider variants of partial I-RAVEN-Mesh dataset with a 468 fraction $q \in \{\frac{1}{64}, \dots, 1\}$ of the training set and compare the performance of a model trained from scratch on a partial dataset to that of a model pre-trained on full I-RAVEN and fine-tuned on a part of I-RAVEN-Mesh. Fig. 6 shows that for $q = \frac{1}{64}$ pre-training RelBase and MRNet on I-RAVEN leads 469 470 to gains smaller than 15 p.p., whereas pre-training DRNet, PoNG, SCL and PredRNet improved 471 their accuracy by 50.6, 19.6, 34.8 and 23.4 p.p., resp. In addition, TL clearly improved performance 472 of DRNet, SCL and PredRNet in all considered settings, in particular for q = 1 by 9.0, 5.9 and 20.5 473 p.p., resp., indicating their capacity for knowledge reuse. 474

475 **Ablation study.** We performed an ablation study with simplified PoNG variants. Table 3 presents 476 the results. The removal of P1 and P2 leads to performance drop, in particular on I-RAVEN-Mesh 477 (-14.9 p.p.) and A/Size (-15.2 p.p.). Similarly, removing P3 and P4 reduces model performance, 478 especially on A/Type (-5.5 p.p.). Disabling TCN leads to generally worse results, primarily on 479 A/Color (-4.9 p.p.) and A/Size (-6.9 p.p.). Training without \mathcal{P}_2^r ($\gamma = 0$) or \mathcal{P}_1^r ($\beta = 0$) 480 typically reduces model performance, but training with one of these rule-based prediction heads 481 compensates to some degree the lack of the other. However, the removal of both ($\gamma = 0 \land \beta = 0$) 482 deteriorates results across all datasets, signifying high relevance of the auxiliary training signal in 483 PoNG's training. To confirm the inherent out-of-distribution generalization abilities of PoNG, we evaluated the model on all PGM regimes without performing any hyperparameter optimization (we 484 only changed the batch size to 256 to reduce training time). Table 4 shows that PoNG achieves strong 485 results on PGM, particularly on the Held-out Triple Pairs regime, exceeding the best reference model Table 5: Visual Analogy Problems (Hill et al., 2019). Results of LBC, NSM, and PredRNet come
from (Yang et al., 2023a, Table 2d). For PoNG, we present mean and std of test accuracy for three
random seeds. ND denotes Novel Domain, NTD — Novel Target Domain, NAV — Novel Attribute
Values, Inter. — Interpolation, Extra. — Extrapolation.

	ND Transfer	NTD LineType	NTD ShapeColor	NAV Inter.	NAV Extra.	Avg
LBC	0.87 ± 0.005	0.76 ± 0.020	0.78 ± 0.004	0.93 ± 0.004	0.62 ± 0.020	0.79
NSM	0.88	0.79	0.78	0.93	0.74	0.82
PredRNet	0.96 ± 0.003	$\textbf{0.82} \pm 0.010$	0.80 ± 0.010	0.97 ± 0.002	0.72 ± 0.060	0.85
PoNG (ours)	$\textbf{0.98} \pm 0.001$	0.78 ± 0.006	$\textbf{0.81} \pm 0.006$	$\textbf{0.98} \pm 0.000$	0.68 ± 0.007	<u>0.84</u>

Table 6: VASR (Bitton et al., 2023). Results of selected baselines come from (Bitton et al., 2023, Table 3). For PoNG, we present mean with std and best-of-3 test accuracy for three random seeds. Sup. denotes Supervised.

Distractors	Zero-Shot ViT	Zero-Shot Swin	Sup. Concat	PoNG (best-of-3)	PoNG (mean \pm std)
Random Difficult	86.0 50.3		$\begin{array}{c} 84.1 \\ 54.9 \end{array}$	92.0 70.5	$91.8 \pm 0.3 \\ 69.5 \pm 1.1$

by 19.4 p.p. We conclude that strong performance of PoNG on A-I-RAVEN and I-RAVEN-Mesh should not be attributed to any specific bias of the model towards these two datasets.

508 Synthetic visual analogies. The VAP benchmark (Hill et al., 2019) was introduced to assess 509 the analogy-based reasoning capabilities of the learning systems. It comprises five generalization 510 regimes: Novel Domain Transfer, Novel Target Domain: Colour of Shapes, Novel Target Domain: 511 Type of Lines, Novel Attribute Values: Interpolation, Novel Attribute Values: Extrapolation, which 512 test the model's generalization to novel domains or attribute values. We compared PoNG to re-513 sults reported in (Yang et al., 2023a, Table 2d), a recent paper introducing the PredRNet model 514 that achieves SOTA results across most VAP regimes. We run PoNG with three random seeds and 515 present its average test accuracy and standard deviation. The results are showcased in Table 5. PoNG presents best results in 3 out of 5 settings, showing its applicability to AVR tasks beyond RPMs. 516

517

497

498

499 500 501

504 505

506

507

Real-world visual analogies. The VASR dataset (Bitton et al., 2023) presents visual analogies 518 comprising real-world images. In effect, the learner needs to additionally understand a rich real-519 world scene, before attempting to solve the presented analogy problem. Following the approach 520 proposed by the VASR authors, we employed the Vision Transformer (ViT) (Dosovitskiy et al., 521 2021) as a perception backbone that produces image embeddings. Specifically, we used the same 522 model variant as (Bitton et al., 2023), which is ViT-L/32 pre-trained on ImageNet-21k at resolution 523 224x224 and fine-tuned on ImageNet-1k at resolution 384x384. We replaced the panel encoder of 524 PoNG with this frozen pre-trained backbone and trained the rest of the model from scratch. We 525 evaluated the model on two VASR splits including random and difficult distractors, resp. As shown 526 in Table 6, in both cases our model outcompetes the strongest result among baselines with 92.0%527 vs. 86.0% and 70.5% vs. 54.9%, resp. The results support the claim that PoNG is a versatile model 528 with strong analogical reasoning capabilities, applicable to both synthetic and real-world domains.

529 530

531

5 CONCLUSION

We investigate generalization capabilities of DL models in the AVR domain. To accelerate research in this area, we propose two RPM benchmarks. Attributeless-I-RAVEN introduces 10 generalization regimes of variable complexity that assess model's capability to solve matrices with rules applied to novel attributes. I-RAVEN-Mesh overlays line-based patterns on top of the RPM, facilitating TL studies. Experiments on 13 strong literature AVR models reveal their limitations in terms of generalization. To elevate state-of-the-art, we introduce PoNG, a novel AVR model capitalizing on parallel design, weight sharing, and normalization. PoNG outcompetes all baselines on the presented challenges, and achieves significant improvement over SOTA reference models on PGM. Furthermore, PoNG excels in solving visual analogy problems comprising both synthetic and real-world images.

540 REFERENCES 541

547

553

554

555

559

565

566

567

569

- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. arXiv preprint 542 arXiv:1607.06450, 2016. 543
- 544 David Barrett, Felix Hill, Adam Santoro, Ari Morcos, and Timothy Lillicrap. Measuring abstract reasoning in neural networks. In International Conference on Machine Learning, pp. 511–520. 546 PMLR, 2018.
- Yaniv Benny, Niv Pekar, and Lior Wolf. Scale-localized abstract reasoning. In Proceedings of the 548 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12557–12565, 2021. 549
- 550 Yonatan Bitton, Ron Yosef, Eliyahu Strugo, Dafna Shahaf, Roy Schwartz, and Gabriel Stanovsky. 551 Vasr: Visual analogies of situation recognition. In Proceedings of the AAAI Conference on Artifi-552 cial Intelligence, volume 37, pp. 241-249, 2023.
 - Mikhail Moiseevich Bongard. The recognition problem. Technical report, Foreign Technology Div Wright-Patterson AFB Ohio, 1968.
- 556 Patricia A Carpenter, Marcel A Just, and Peter Shell. What one intelligence test measures: a theoretical account of the processing in the raven progressive matrices test. Psychological review, 97 558 (3):404, 1990.
- François Chollet. On the measure of intelligence. arXiv preprint arXiv:1911.01547, 2019. 560
- 561 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 562 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszko-563 reit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In International Conference on Learning Representations, 2021.
- SM Eslami, Nicolas Heess, Theophane Weber, Yuval Tassa, David Szepesvari, Geoffrey E Hinton, et al. Attend, infer, repeat: Fast scene understanding with generative models. Advances in neural information processing systems, 29, 2016. 568
 - Thomas G Evans. A heuristic program to solve geometric-analogy problems. In Proc. of the April 21-23, 1964, spring joint computer conference, pp. 327–338, 1964.
- 571 François Fleuret, Ting Li, Charles Dubout, Emma K Wampler, Steven Yantis, and Donald Geman. 572 Comparing machines and humans on a visual categorization test. Proceedings of the National 573 Academy of Sciences, 108(43):17621-17625, 2011. 574
- 575 Harry E Foundalis. Phaeaco: A cognitive architecture inspired by Bongard's problems. PhD dissertation, Indiana University, 2006. 576
- 577 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, 578 Hal Daumé Iii, and Kate Crawford. Datasheets for datasets. Communications of the ACM, 64 579 (12):86-92, 2021. 580
- Dedre Gentner. The structure of analogical models in science. Bolt Beranek and Newman Cam-581 bridge, 1980. 582
- 583 José Hernández-Orallo. The measure of all minds: evaluating natural and artificial intelligence. 584 Cambridge University Press, 2017. 585
- José Hernández-Orallo, Fernando Martínez-Plumed, Ute Schmid, Michael Siebers, and David L 586 Dowe. Computer models solving intelligence test problems: Progress and implications. Artificial Intelligence, 230:74-107, 2016. 588
- 589 Felix Hill, Adam Santoro, David Barrett, Ari Morcos, and Timothy Lillicrap. Learning to make 590 analogies by contrasting abstract relational structure. In International Conference on Learning Representations, 2019. 592
- Dokhyam Hoshen and Michael Werman. IQ of neural networks. arXiv preprint arXiv:1710.01692, 2017.

594 595 596	Sheng Hu, Yuqing Ma, Xianglong Liu, Yanlu Wei, and Shihao Bai. Stratified rule-aware network for abstract visual reasoning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 35, pp. 1567–1574, 2021.
598 599 600	Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In <i>International Conference on Machine Learning</i> , pp. 448–456. PMLR, 2015.
601 602 603	Yifan Jiang, Jiarui Zhang, Kexuan Sun, Zhivar Sourati, Kian Ahrabian, Kaixin Ma, Filip Ilievski, and Jay Pujara. Marvel: Multidimensional abstraction and reasoning through visual evaluation and learning. <i>arXiv preprint arXiv:2404.13591</i> , 2024.
605 606	Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. <i>International Conference on Learning Representations</i> , 2014.
607 608 609	Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. <i>Advances in neural information processing systems</i> , 25, 2012.
610 611 612	Maithilee Kunda, Keith McGreggor, and Ashok Goel. Taking a look (literally!) at the raven's intelligence test: Two visual solution strategies. In <i>Proc. of the Annual Meeting of the Cognitive Science Society</i> , volume 32, 2010.
613 614 615 616	Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. <i>Advances in Neural Information Processing Systems</i> , 33:11525–11538, 2020.
617 618 619	Andrew Lovett, Kenneth Forbus, and Jeffrey Usher. Analogy with qualitative spatial representations can simulate solving raven's progressive matrices. In <i>Proc. of the Annual Meeting of the Cognitive Science Society</i> , volume 29, 2007.
620 621 622	Mikołaj Małkiński and Jacek Mańdziuk. Deep learning methods for abstract visual reasoning: A survey on raven's progressive matrices. <i>arXiv preprint arXiv:2201.12382</i> , 2022.
623 624	Mikołaj Małkiński and Jacek Mańdziuk. A review of emerging research directions in abstract visual reasoning. <i>Information Fusion</i> , 91:713–736, 2023.
626 627	Mikołaj Małkiński and Jacek Mańdziuk. Multi-label contrastive learning for abstract visual reason- ing. <i>IEEE Transactions on Neural Networks and Learning Systems</i> , 35(2):1941–1953, 2024a.
628 629 630	Mikołaj Małkiński and Jacek Mańdziuk. One self-configurable model to solve many abstract visual reasoning problems. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 14297–14305, 2024b.
632 633	Jacek Mańdziuk and Adam Żychowski. DeepIQ: A human-inspired AI system for solving IQ test problems. In 2019 International Joint Conference on Neural Networks, pp. 1–8. IEEE, 2019.
634 635 636 637	Laura E Matzen, Zachary O Benz, Kevin R Dixon, Jamie Posey, James K Kroger, and Ann E Speed. Recreating raven's: Software for systematically generating large numbers of raven-like matrix problems with normed properties. <i>Behavior research methods</i> , 42(2):525–541, 2010.
638 639	Melanie Mitchell. Abstraction and analogy-making in artificial intelligence. <i>Annals of the New York Academy of Sciences</i> , 1505(1):79–101, 2021.
640 641 642	Melanie Mitchell, Alessandro B Palmarini, and Arseny Moskvichev. Comparing humans, gpt-4, and gpt-4v on abstraction and reasoning tasks. <i>arXiv preprint arXiv:2311.09247</i> , 2023.
643 644 645	Shanka Subhra Mondal, Taylor Whittington Webb, and Jonathan Cohen. Learning to reason over visual objects. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=uR6x8Be7o_M.
040 647	Shanka Subhra Mondal, Jonathan D Cohen, and Taylor W Webb. Slot abstractors: Toward scalable abstract visual reasoning. <i>arXiv preprint arXiv:2403.03458</i> , 2024.

- 648 Arsenii Kirillovich Moskvichev, Victor Vikram Odouard, and Melanie Mitchell. The conceptARC 649 benchmark: Evaluating understanding and generalization in the ARC domain. Transactions on 650 Machine Learning Research, 2023. ISSN 2835-8856. 651 Weili Nie, Zhiding Yu, Lei Mao, Ankit B Patel, Yuke Zhu, and Anima Anandkumar. Bongard-652 logo: A new benchmark for human-level concept learning and reasoning. Advances in Neural 653 Information Processing Systems, 33:16468–16480, 2020. 654 Victor Vikram Odouard and Melanie Mitchell. Evaluating understanding on conceptual abstraction 655 benchmarks. arXiv preprint arXiv:2206.14187, 2022. 656 657 Yonggang Qi, Kai Zhang, Aneeshan Sain, and Yi-Zhe Song. Pqa: Perceptual question answering. 658 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 659 12056-12064, 2021. 660 James C Raven. Mental tests used in genetic studies: The performance of related individuals on 661 tests mainly educative and mainly reproductive. Master's thesis, University of London, 1936. 662 663 John C Raven and John Hugh Court. Raven's progressive matrices and vocabulary scales. Oxford 664 pyschologists Press Oxford, England, 1998. 665 Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter 666 Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. Ad-667 vances in neural information processing systems, 30:4967–4976, 2017. 668 669 Murray Shanahan, Kyriacos Nikiforou, Antonia Creswell, Christos Kaplanis, David Barrett, and 670 Marta Garnelo. An explicitly relational neural network architecture. In *International Conference* on Machine Learning, pp. 8593-8603. PMLR, 2020. 671 672 Richard E Snow, Patrick C Kyllonen, and Brachia Marshalek. The topography of ability and learning 673 correlations. Advances in the psychology of human intelligence, 2(S 47):103, 1984. 674 Steven Spratley, Krista Ehinger, and Tim Miller. A closer look at generalisation in raven. In Com-675 puter Vision-ECCV 2020: 16th European Conference, Glasgow, UK, August 23-28, 2020, Pro-676 ceedings, Part XXVII 16, pp. 601-616. Springer, 2020. 677 678 Sebastian Stabinger, David Peer, Justus Piater, and Antonio Rodríguez-Sánchez. Evaluating the 679 progress of deep learning for visual relational concepts. Journal of Vision, 21(11):8–8, 2021. 680 Claes Strannegård, Simone Cirillo, and Victor Ström. An anthropomorphic method for progressive 681 matrix problems. Cognitive Systems Research, 22:35–46, 2013. 682 683 Paulina Tomaszewska, Adam Żychowski, and Jacek Mańdziuk. Duel-based deep learning system 684 for solving iq tests. In International Conference on Artificial Intelligence and Statistics, pp. 10483-10492. PMLR, 2022. 685 686 Han LJ van der Maas, Lukas Snoek, and Claire E Stevenson. How much intelligence is there in 687 artificial intelligence? a 2020 update. Intelligence, 87:101548, 2021. 688 Ke Wang and Zhendong Su. Automatic generation of raven's progressive matrices. In Twenty-fourth 689 international joint conference on artificial intelligence, 2015. 690 691 Taylor Webb, Zachary Dulberg, Steven Frankland, Alexander Petrov, Randall O'Reilly, and 692 Jonathan Cohen. Learning representations that support extrapolation. In International Conference 693 on Machine Learning, pp. 10136–10146. PMLR, 2020. 694 Yuhuai Wu, Honghua Dong, Roger Grosse, and Jimmy Ba. The scattering compositional learner: Discovering objects, attributes, relationships in analogical reasoning. arXiv preprint 696 arXiv:2007.04212, 2020. 697 Lingxiao Yang, Hongzhi You, Zonglei Zhen, Dahui Wang, Xiaohong Wan, Xiaohua Xie, and Ru-Yuan Zhang. Neural prediction errors enable analogical visual reasoning in human standard intel-699 ligence tests. In Proceedings of the 40th International Conference on Machine Learning, volume 700 202 of Proceedings of Machine Learning Research, pp. 39572–39583. PMLR, 23–29 Jul 2023a.
 - URL https://proceedings.mlr.press/v202/yang23r.html.

Yuan Yang, Deepayan Sanyal, Joel Michelson, James Ainooson, and Maithilee Kunda. A conceptual chronicle of solving raven's progressive matrices computationally. In Proceedings of the 8th International Workshop on Artificial Intelligence and Cognition, 2022. Yuan Yang, Deepayan Sanyal, James Ainooson, Joel Michelson, Effat Farhana, and Maithilee Kunda. A cognitively-inspired neural architecture for visual abstract reasoning using contrastive perceptual and conceptual processing. arXiv preprint arXiv:2309.10532, 2023b. Chi Zhang, Feng Gao, Baoxiong Jia, Yixin Zhu, and Song-Chun Zhu. RAVEN: A dataset for rela-tional and analogical visual reasoning. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5317–5327, 2019a. Chi Zhang, Baoxiong Jia, Feng Gao, Yixin Zhu, Hongjing Lu, and Song-Chun Zhu. Learning perceptual inference by contrasting. Advances in neural information processing systems, 32: 1075-1087, 2019b. Chi Zhang, Baoxiong Jia, Song-Chun Zhu, and Yixin Zhu. Abstract spatial-temporal reasoning via probabilistic abduction and execution. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 9736–9746, 2021. Chi Zhang, Sirui Xie, Baoxiong Jia, Ying Nian Wu, Song-Chun Zhu, and Yixin Zhu. Learning alge-braic representation for systematic generalization in abstract reasoning. In European Conference on Computer Vision, pp. 692-709. Springer, 2022a. Wenbo Zhang, Site Mo, Xianggen Liu, Sen Song, et al. Learning robust rule representations for abstract reasoning via internal inferences. Advances in Neural Information Processing Systems, 35:33550-33562, 2022b. Wenhe Zhang, Chi Zhang, Yixin Zhu, and Song-Chun Zhu. Machine number sense: A dataset of visual arithmetic problems for abstract and relational reasoning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pp. 1332–1340, 2020. Kai Zhao, Chang Xu, and Bailu Si. Learning visual abstract reasoning through dual-stream net-works. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 38, pp. 16979– 16988, 2024. Tao Zhuo and Mohan Kankanhalli. Effective abstract reasoning with dual-contrast network. In International Conference on Learning Representations, 2021.



Figure 7: Attributeless-I-RAVEN. Left: Matrices from the A/Type regime belonging to the Center configuration. In (a), object type is constant across rows, while in (b) it's governed by the Distribute Three rule. Right: Matrices from the A/Size regime belonging to the Out-InCenter configuration. In (c), object size is constant across rows in both inner and outer image parts, while in (d) the inner and outer components are governed by the Arithmetic and Progression rules, resp.



Figure 8: Attributeless-I-RAVEN. Left: Matrices from the A/ColorSize regime belonging to the Up-Down configuration. In (a), object color and size is constant across rows in both components, while in (b) they are governed by Progression and Distribute Three in the upper component, resp., and by Distribute Three in the lower one. Right: Matrices from the A/ColorType regime belonging to the 2x2 Grid configuration. In (c), object color and type is constant across rows, while in (d) they are governed by the Distribute Three rule.

A ADDITIONAL MATRIX EXAMPLES

796 797

798

799

800

801

794

769

770

771

772

773

774

Figure 7 presents matrix examples from A/Type and A/Size, the primary regimes of Attributeless-I-RAVEN. Figures 8, 9 and 10 depict matrix examples from the extended regimes of Attributeless-I-RAVEN: A/ColorSize and A/ColorType (Fig. 8), A/SizeType and A/Color-Progression (Fig. 9), and A/Color-Arithmetic and A/Color-DistributeThree (Fig. 10). Figure 11 presents matrix examples from I-RAVEN-Mesh concerning the Number attribute.

802 803 804

805

B LIMITATIONS AND FUTURE WORK

In this work we study generalization and knowledge transfer in contemporary AVR models employ ing RPM datasets, and compare the introduced PoNG model with SOTA models in solving visual
 analogy problems. However, the set of problems in the AVR domain also includes other tasks not
 covered in the paper (Małkiński & Mańdziuk, 2023). The Machine Number Sense dataset presents
 visual arithmetic problems (Zhang et al., 2020), VAEC defines an extrapolation challenge (Webb



Figure 9: Attributeless-I-RAVEN. Left: Matrices from the A/SizeType regime belonging to the Center configuration. In (a), object size and type are constant across rows, while in (b) they are governed by the Progression rule. Right: Matrices from the A/Color-Progression regime belonging to the 3x3 Grid configuration. In (c), object color is governed by the Progression rule, while in (d) by Distribute Three.



Figure 10: Attributeless-I-RAVEN. Left: Matrices from the A/Color-Arithmetic regime belonging to the Out-InCenter configuration. In (a), object color in the inner component is governed by Arithmetic, while in (b) it is governed by Distribute Three. Right: Matrices from the A/Color-DistributeThree regime belonging to the Left-Right configuration. In (c), object color is governed by the Distribute Three rule in both components, while in (d) by Constant in the left component and by Arithmetic in the right one.

et al., 2020), while ARC proposes a set of diverse tasks in a few-shot learning setting (Chollet, 2019). Future research in this area may juxtapose the performance of AVR models across a set of benchmarks oriented towards generalization to ensure generalization advances beyond the RPM and visual analogy datasets.

In the paper we claim that the proposed pathways block, a key component of the introduced model, is a generic module also applicable to other tasks that require reasoning over a set of objects (vector embeddings). Nevertheless, the experimental evaluation of PoNG presented in the paper is focused on RPM benchmarks, including I-RAVEN, I-RAVEN-Mesh, Attributeless-I-RAVEN, and PGM, and two visual analogy datasets, i.e. VAP and VASR. Assessing model's performance on other problems constitutes an interesting extension of this work.

C DATASET DETAILS

Rule encoding. As discussed in Section 3.3, we use sparse encoding (Małkiński & Mańdziuk, 2024a) to represent the set of matrix rules as a vector $r \in \mathbb{R}^{d_r}$, such that $d_r = 48$ for I-RAVEN-Mesh and $d_r = 40$ otherwise. The set of rules \mathcal{R} in I-RAVEN is {Constant, Progression, Arithmetic, Distribute Three} and the set of attributes \mathcal{A} is {Position, Number,

875 876

877

878

879

880

881 882 883

884

885

886

887

889

890

891

892

893 894

895 896

897 898 899

900 901



Figure 11: **I-RAVEN-Mesh.** The examples showcase matrices with the Number attribute of the mesh component governed by all applicable rules. (a) Line number is constant in each row. (b) The number of lines increases by 2 from left to right. (c) The number of lines in the third column is the difference between the number of lines in the second and first columns. (d) The numbers of lines in each row compose a set $\{2, 3, 6\}$.

Type, Size, Color}. It follows that there is $|\mathcal{R}| \times |\mathcal{A}| = 20$ unique rule-attribute pairs. In addition, the Left-Right, Up-Down, Out-InCenter, and Out-InGrid configurations in I-RAVEN comprise two components in which rules exist independently, e.g., the Left-Right component contains matrices with separate rules applied to the left and right sides. This gives an upper bound of 40 rule-attribute combinations in each configuration. As discussed in Section 3.2 and presented in Table 1, the Mesh component introduced in I-RAVEN-Mesh comprises two attributes and four rules, leading to a total of 48 rule-attribute combinations per configuration. As an example, in the Up-Down configuration of I-RAVEN-Mesh, there are 20 rule-attribute combinations for the upper component, another 20 for the lower component, and 8 for the Mesh component. The sparse encoding encodes each rule in a matrix as a one-hot vector and applies the OR operation to the set of one-hot vectors, producing a multi-hot representation of matrix rules.

D MODEL DETAILS

Tables 7, 8, and 9 list all PoNG hyperparameters.

E EXTENDED RESULTS

Fig. 12 presents extended results of Fig. 6.

Table 10 shows the aggregated performance of all considered models on 6 extended A-I-RAVEN regimes.

Tables 11 - 22 present the results (mean and standard deviation) of all considered models on test 905 and validation splits and the difference between these two splits for particular datasets/regimes. The 906 results support the analysis of dataset difficulty presented in Section 4. The difference in model per-907 formance between test and validation splits in I-RAVEN (Table 11) and I-RAVEN-Mesh (Table 12) 908 is negligible. In Attributeless-I-RAVEN regimes, however, the difference is significant, showing 909 limitations of all evaluated models in terms of generalization. Across 4 primary regimes (Tables 13 – 910 16), the biggest difference concerns the A/Type regime, suggesting that generalization of rules ap-911 plied to novel shape types constitutes a real challenge for the contemporary models. In all 3 extended 912 regimes concerning held-out attribute pairs (A/ColorSize, A/ColorType, and A/SizeType) 913 the performance difference on test and validation splits is bigger than in the primary regimes (see 914 Tables 17 - 19). This drop stems from overall weaker performance on the test split, confirming high 915 difficulty of these regimes. Model performance on the next 3 regimes concerning the Color attribute and rules other than Constant (A/Color-Progression, A/Color-Arithmetic, 916 and A/Color-DistributeThree) is better, though further progress in generalization is re-917 quired to fully close the performance gap between test and validation splits (see Tables 20 - 22).

Table 7: **PoNG hyperparameters: Panel encoder** \mathcal{E} . The parameters of convolution layers are denoted as [# input channels \rightarrow # output channels, kernel size, stride, padding]; of pooling layers as [kernel size, stride, padding]; of linear layers as [# input neurons \rightarrow # output neurons]; of flatten operators as [input dimensions \rightarrow output dimensions]; of position embedding as [dimensionality of the position embedding vector].

923		
924	LAYER	Hyperparameters
925	CONV2D-RELU-BN2D	$[1 \rightarrow 32, 7 \times 7, 2 \times 2, 3 \times 3]$
926	CONV2D-RELU-BN2D	$[32 \rightarrow 32, 7 \times 7, 2 \times 2, 3 \times 3]$
927	MaxPool	$[3 \times 3, 2 \times 2, 1 \times 1]$
928	MAXPOOL	$\begin{bmatrix} 3 \times 3, 2 \times 2, 1 \times 1 \end{bmatrix}$
929	CONV2D Sum	$[1 \rightarrow 32, 1 \times 1, 1 \times 1, 0 \times 0]$
930		
931	CONV2D-RELU-BN2D CONV2D-RELU-BN2D	$\begin{bmatrix} 32 \rightarrow 32, 7 \times 7, 2 \times 2, 3 \times 3 \\ 32 \rightarrow 32, 7 \times 7, 2 \times 2, 3 \times 3 \end{bmatrix}$
932	MAXPOOL	$\begin{bmatrix} 52 & 7 & 52, 7 & 7, 2 & 2, 5 & 7 & 6 \\ & [3 \times 3, 2 \times 2, 1 \times 1] \end{bmatrix}$
933	MAXPOOL	$[3 \times 3, 2 \times 2, 1 \times 1]$
934	Conv2D	$[32 \rightarrow 32, 1 \times 1, 1 \times 1, 0 \times 0]$
935	SUM	
936	FLATTEN (HEIGHT & WIDTH)	$[5 \times 5 \to 25]$
937	LINEAR-RELU	$[25 \rightarrow 25]$
938	FLATTEN (DEPTH & SPATIAL)	$[32 \times 25 \rightarrow 800]$
939	LN LNEAD DELU	
940	LINEAR-KELU I N	$[800 \rightarrow 1000]$
941	LINEAR	$[1600 \rightarrow 800]$
942	Sum	L]
943	POSITION EMBEDDING	[25]

944 945

~~~

Tables 23 - 34 present the results (mean and standard deviation) of all considered models in detail 946 for all matrix configurations. The most challenging configurations in I-RAVEN and I-RAVEN-947 Mesh are 3x3Grid and Out-InGrid, in which image panels contain more objects than in 948 the remaining configurations. Apparently, such setups require stronger reasoning capabilities to 949 correctly identify the rules applied to multiple objects. Also, the results on the Left-Right 950 and Up-Down configurations are relatively weaker in majority of regimes. In these configura-951 tions, rules may be applied to both matrix components (left/right and up/down, resp.), which de-952 mands stronger reasoning capabilities. This also concerns the Out-InGrid configuration in the 953 A/Size regime, and the Out-InCenter configuration in the A/SizeType regime. Results in the A/Position regime are close-to-perfect in configurations comprising a single object in 954 each component (Center, Left-Right, Up-Down, and Out-InCenter) and weaker in the 955 remaining configurations (2x2Grid, 3x3Grid and Out-InGrid). This performance drop can 956 be attributed to the fact that Position attribute can only be effectively applied to the 2x2Grid, 957 3x3Grid and Out-InGrid configurations allowing modification of the object's position. In the 958 remaining configurations its application does not introduce any changes. 959

- 961
- 962 963
- 964
- 965 966
- 967
- 968
- 969
- 970
- 971

Table 8: **PoNG hyperparameters: Reasoner**  $\mathcal{R}$ **.** The parameters of convolution layers are denoted as [# input channels  $\rightarrow$  # output channels, kernel size, stride, padding]; of group and group-pair convolution layers as [# input channels  $\rightarrow$  # output channels, kernel size, stride, padding, groups]; of pooling layers as [kernel size, stride, padding]; of linear layers as [# input neurons  $\rightarrow$  # output neurons]; of flatten operators as [input dimensions  $\rightarrow$  output dimensions].

| 979  |                               |                                                                                                |
|------|-------------------------------|------------------------------------------------------------------------------------------------|
| 980  | LAYER                         | HYPERPARAMETERS                                                                                |
| 981  | STACK                         |                                                                                                |
| 082  | BN1D                          |                                                                                                |
| 002  | CONV1D                        | $[9 \rightarrow 32, 1, 1, 0, \text{bias} = \text{False}]$                                      |
| 903  | CONV1D-RELU-BN1D              | $[9 \to 32, 7, 1, 3]$                                                                          |
| 984  | G-C-TCN-RELU-BN1D             | $[3 \rightarrow 32, 7, 1, 3, 3]$                                                               |
| 985  | GP-C-TCN-RELU-BN1D            | $[6 \rightarrow 32, 7, 1, 3, 3]$                                                               |
| 986  | AVGPOOL1D                     | [10, 8, 1]                                                                                     |
| 987  | BN1D                          |                                                                                                |
| 988  | Conv1D                        | $[32 \rightarrow 32, 1, 1, 0, \text{bias} = \text{False}]$                                     |
| 989  | CONV1D-RELU-BN1D              | $[32 \rightarrow 32, 7, 1, 3]$                                                                 |
| 990  | CONV1D-RELU-BN1D              | $[32 \rightarrow 32, 7, 1, 3]$                                                                 |
| 001  | G-C-TCN-RELU-BN1D             | $[4 \rightarrow 32, 7, 1, 3, 8]$                                                               |
| 331  | G-C-TCN-RELU-BN1D             | $[4 \rightarrow 32, 7, 1, 3, 8]$                                                               |
| 992  | GP-C-TCN-RELU-BN1D            | $[16 \rightarrow 32, 7, 1, 3, 4]$                                                              |
| 993  | GP-C-ICN-RELU-BNID            | $[16 \to 32, 7, 1, 3, 4]$                                                                      |
| 994  | AvgPool1D                     | [6, 4, 1]                                                                                      |
| 995  | BN1D                          |                                                                                                |
| 996  | Conv1D                        | $[32 \rightarrow 32, 1, 1, 0, \text{bias} = \text{False}]$                                     |
| 997  | CONV1D-RELU-BN1D              | $[32 \rightarrow 32, 7, 1, 3]$                                                                 |
| 998  | CONV1D-RELU-BN1D              | $[32 \rightarrow 32, 7, 1, 3]$                                                                 |
| 000  | G-C-TCN-RELU-BN1D             | $[4 \rightarrow 32, 7, 1, 3, 8]$                                                               |
| 1000 | G-C-ICN-KELU-BNID             | $[4 \rightarrow 32, 7, 1, 3, 8]$                                                               |
| 1000 | GP-C-ICN-KELU-BNID            | $\begin{bmatrix} 10 \rightarrow 32, 7, 1, 3, 4 \\ 16 \rightarrow 22, 7, 1, 2, 4 \end{bmatrix}$ |
| 1001 | GP-C-ICN-RELU-DNID            | $[10 \to 52, 7, 1, 5, 4]$                                                                      |
| 1002 | ADAPTIVE AVGPOOL1D            | $[25 \rightarrow 16]$                                                                          |
| 1003 | FLATTEN (DEPTH & FEATURE DIM) | $[32 \times 16 \rightarrow 512]$                                                               |
| 1004 | LINEAR-RELU-BN1D              | $\begin{bmatrix} 512 \rightarrow 512 \end{bmatrix}$                                            |
| 1005 | LINEAR                        | $[512 \rightarrow 128]$                                                                        |

Table 9: **PoNG hyperparameters: Prediction heads.** The parameters of linear layers are denoted as [# input neurons  $\rightarrow$  # output neurons].

| LAYER                       | <b>HYPERPARAMETERS</b>                                 |
|-----------------------------|--------------------------------------------------------|
| Target head $\mathcal{P}^y$ |                                                        |
| LINEAR-RELU-LINE<br>LINEAR  | $\begin{array}{c c c c c c c c c c c c c c c c c c c $ |
| AGGREGATE RULE HE           | EAD $\mathcal{P}_1^r$                                  |
| Sum                         |                                                        |
| LINEAR-RELU                 | $[128 \rightarrow 128]$                                |
| LINEAR                      | $[128 \rightarrow d_r]$                                |
| TARGET-CONDITIONE           | ED RULE HEAD $\mathcal{P}_2^r$                         |
| WEIGHTED SUM                |                                                        |
| LINEAR                      | $[128 \rightarrow d_r]$                                |

Table 10: A-I-RAVEN extended regimes. CS, CT, and ST denote ColorSize, ColorType, and SizeType, resp. P, A, and D3 denote Progression, Arithmetic, and Distribute Three, resp.

|             | A/CS                      | A/CT                      | A/ST                    | A/Color-P               | A/Color-A           | A/Color-D3         |
|-------------|---------------------------|---------------------------|-------------------------|-------------------------|---------------------|--------------------|
| ALANS       | $  15.1 (\pm 3.3)$        | $17.7 (\pm 3.2)$          | $15.7 (\pm 3.2)$        | $24.8 (\pm 18.8)$       | $18.3 (\pm 6.6)$    | $22.4 (\pm 7.7)$   |
| CPCNet      | $33.0(\pm 5.3)$           | $25.0 (\pm 0.9)$          | $24.1 (\pm 1.2)$        | $50.5 (\pm 0.6)$        | $45.9 (\pm 2.7)$    | $37.8 (\pm 0.9)$   |
| CNN-LSTM    | $13.4 (\pm 0.9)$          | $14.7 (\pm 1.7)$          | $13.0 (\pm 0.1)$        | $17.2 (\pm 1.5)$        | $17.1 (\pm 3.7)$    | $20.6 (\pm 6.7)$   |
| CoPINet     | $18.3 (\pm 0.3)$          | $17.2 (\pm 0.1)$          | $19.7 (\pm 0.7)$        | $35.8 (\pm 0.6)$        | $35.2 (\pm 0.5)$    | $26.9 (\pm 0.5)$   |
| DRNet       | $38.3 (\pm 0.5)$          | $29.5 (\pm 0.5)$          | $31.6 (\pm 1.2)$        | $72.8 (\pm 1.3)$        | $66.7 (\pm 1.2)$    | $63.2 (\pm 0.3)$   |
| MRNet       | $18.7 (\pm 1.1)$          | $20.0 (\pm 2.6)$          | $28.2 (\pm 0.9)$        | $34.4 (\pm 3.4)$        | $35.7 (\pm 5.9)$    | $18.6 (\pm 0.1)$   |
| PrAE        | $30.0(\pm 1.1)$           | $26.7 (\pm 0.7)$          | $25.6 (\pm 0.8)$        | $62.3 (\pm 0.9)$        | $43.0 (\pm 26.5)$   | $55.1 (\pm 0.8)$   |
| PredRNet    | $31.0(\pm 1.6)$           | $28.0 (\pm 0.7)$          | $27.9 (\pm 0.5)$        | $62.3 (\pm 2.2)$        | $56.9(\pm 1.4)$     | $48.5 (\pm 0.9)$   |
| RelBase     | $36.6 (\pm 0.8)$          | $29.7 (\pm 0.6)$          | $31.1 (\pm 1.0)$        | $73.0 (\pm 1.8)$        | $66.2 (\pm 1.0)$    | $65.7 (\pm 4.6)$   |
| SCL         | $40.8(\pm 3.2)$           | $32.0(\pm 2.3)$           | <b>33.5</b> $(\pm 0.7)$ | $75.6(\pm 10.1)$        | $60.0 (\pm 4.1)$    | $63.9(\pm 4.3)$    |
| SRAN        | $\overline{22.7}$ (± 1.1) | $\overline{20.9}$ (± 0.9) | $23.3 (\pm 0.3)$        | $42.1 (\pm 2.3)$        | $39.9(\pm 2.7)$     | $34.6 (\pm 3.6)$   |
| STSN        | $27.3 (\pm 4.6)$          | $21.9 (\pm 4.6)$          | $12.3 (\pm 0.1)$        | $39.9(\pm 14.7)$        | $25.7 (\pm 10.6)$   | $20.7 (\pm 7.7)$   |
| WReN        | $13.5(\pm 0.1)$           | $13.8 (\pm 0.7)$          | $14.1 (\pm 0.2)$        | $18.0(\pm 0.4)$         | $17.1 (\pm 0.2)$    | $17.7 (\pm 0.6)$   |
| PoNG (ours) | <b>44.7</b> (± 2.1)       | $34.3 (\pm 0.8)$          | $32.1(\pm 2.1)$         | <b>81.4</b> $(\pm 3.1)$ | <b>70.0</b> (± 4.1) | <b>81.3</b> (±1.6) |
|             |                           |                           |                         |                         |                     |                    |

Table 11: I-RAVEN.

Table 12: I-RAVEN-Mesh.

|                 | Test                | Val                     | Test - Val |             | Test                | Val                        | Test - Val |
|-----------------|---------------------|-------------------------|------------|-------------|---------------------|----------------------------|------------|
| ALANS           | $27.0 (\pm 8.4)$    | $27.0 (\pm 8.6)$        | + 0.1      | ALANS       | $15.9(\pm 2.6)$     | $17.1 (\pm 3.6)$           | - 1.3      |
| CPCNet          | $70.4 (\pm 6.4)$    | $69.6 (\pm 6.9)$        | + 0.7      | CPCNet      | $66.6 (\pm 5.1)$    | $66.5 (\pm 5.4)$           | + 0.1      |
| <b>CNN-LSTM</b> | 27.5 $(\pm 1.5)$    | $27.4 (\pm 1.7)$        | + 0.1      | CNN-LSTM    | $28.9 (\pm 0.4)$    | $29.3 (\pm 0.6)$           | -0.4       |
| CoPINet         | $43.2 (\pm 0.1)$    | $42.5 (\pm 0.6)$        | + 0.7      | CoPINet     | $41.1 (\pm 0.3)$    | $41.3 (\pm 0.2)$           | -0.2       |
| DRNet           | $90.9(\pm 1.1)$     | $90.8 (\pm 1.2)$        | + 0.2      | DRNet       | $83.9(\pm 2.7)$     | $84.2 (\pm 2.6)$           | -0.3       |
| MRNet           | $86.7 (\pm 2.3)$    | $86.3 (\pm 2.0)$        | + 0.5      | MRNet       | $79.5 (\pm 2.0)$    | $80.5~(\pm 2.5)$           | -1.0       |
| PrAE            | $19.5 (\pm 0.4)$    | $19.4 (\pm 0.8)$        | + 0.0      | PrAE        | $33.2 (\pm 0.4)$    | $33.0 (\pm 0.9)$           | + 0.1      |
| PredRNet        | $88.8 (\pm 1.8)$    | $88.3 (\pm 1.9)$        | + 0.5      | PredRNet    | $59.2 (\pm 6.4)$    | $59.3 (\pm 6.9)$           | - 0.0      |
| RelBase         | $89.6 (\pm 0.6)$    | $89.5 (\pm 0.5)$        | + 0.1      | RelBase     | $84.9(\pm 4.4)$     | $\underline{85.0}$ (± 4.5) | - 0.1      |
| SCL             | $83.4(\pm 2.5)$     | $83.0 (\pm 2.5)$        | + 0.4      | SCL         | $80.9 (\pm 1.5)$    | $81.0 (\pm 1.5)$           | - 0.1      |
| SRAN            | $58.2 (\pm 1.6)$    | $58.0 (\pm 1.3)$        | + 0.2      | SRAN        | $57.8 (\pm 0.2)$    | $58.0 (\pm 0.3)$           | -0.2       |
| STSN            | $59.0 (\pm 18.5)$   | $59.1 (\pm 18.4)$       | - 0.1      | STSN        | $48.7 (\pm 11.5)$   | $48.8 (\pm 10.9)$          | - 0.1      |
| WReN            | $18.4 (\pm 0.0)$    | $18.5 (\pm 0.3)$        | - 0.1      | WReN        | $25.7 (\pm 0.2)$    | $25.6 (\pm 0.4)$           | + 0.0      |
| PoNG (ours)     | <b>95.9</b> (± 0.7) | <b>95.6</b> $(\pm 0.7)$ | + 0.3      | PoNG (ours) | <b>89.3</b> (± 2.4) | <b>89.1</b> (± 2.5)        | + 0.3      |

Table 13: A/Color.

Table 14: A/Position.

|             | Test                | Val               | Test-Val |             | Test                       | Val                | Test - Val |
|-------------|---------------------|-------------------|----------|-------------|----------------------------|--------------------|------------|
| ALANS       | $ 15.2(\pm 1.4) $   | $16.4 (\pm 2.1)$  | - 1.2    | ALANS       | $16.0 (\pm 1.0)$           | $15.2 (\pm 1.3)$   | + 0.8      |
| CPCNet      | $51.2 (\pm 3.8)$    | $77.0(\pm 6.1)$   | -25.7    | CPCNet      | $68.3 (\pm 4.0)$           | $90.6 (\pm 5.3)$   | -22.3      |
| CNN-LSTM    | $17.0 (\pm 3.1)$    | $31.0(\pm 4.4)$   | -13.9    | CNN-LSTM    | $24.0 (\pm 2.9)$           | $36.4 (\pm 3.7)$   | -12.4      |
| CoPINet     | $32.5 (\pm 0.2)$    | $49.9 (\pm 0.7)$  | -17.4    | CoPINet     | $41.3 (\pm 1.6)$           | $54.7 (\pm 1.7)$   | -13.4      |
| DRNet       | $70.0 (\pm 1.6)$    | $95.4 (\pm 0.2)$  | -25.4    | DRNet       | $\underline{77.5}$ (± 0.9) | $97.8(\pm 0.1)$    | -20.2      |
| MRNet       | $33.6 (\pm 8.2)$    | $86.2 (\pm 6.6)$  | -52.6    | MRNet       | $62.6 (\pm 2.6)$           | $94.4 (\pm 7.0)$   | -31.8      |
| PrAE        | $47.9 (\pm 0.9)$    | $60.9 (\pm 1.4)$  | -13.0    | PrAE        | $68.2 (\pm 3.3)$           | $80.1 (\pm 3.4)$   | -12.0      |
| PredRNet    | $59.4 (\pm 1.0)$    | $92.2 (\pm 1.0)$  | -32.9    | PredRNet    | $73.7 (\pm 0.7)$           | $97.4 \ (\pm 0.5)$ | -23.6      |
| RelBase     | $67.4 (\pm 2.7)$    | $95.2 (\pm 0.4)$  | -27.8    | RelBase     | $76.6 (\pm 0.3)$           | $97.0 \ (\pm 0.2)$ | -20.4      |
| SCL         | $65.1 (\pm 2.0)$    | $84.4 (\pm 0.5)$  | -19.2    | SCL         | $76.7 (\pm 7.1)$           | $94.7 (\pm 5.3)$   | -18.0      |
| SRAN        | $38.3 (\pm 1.0)$    | $63.7 (\pm 0.3)$  | -25.4    | SRAN        | $56.9 (\pm 0.7)$           | $75.6 (\pm 1.4)$   | -18.8      |
| STSN        | $39.3 (\pm 6.9)$    | $71.3 (\pm 17.0)$ | -32.0    | STSN        | $36.1 (\pm 19.9)$          | $50.7~(\pm 27.3)$  | -14.6      |
| WReN        | $16.9 (\pm 0.5)$    | $23.2 (\pm 0.8)$  | -6.3     | WReN        | $17.3 (\pm 0.4)$           | $23.3 (\pm 0.5)$   | - 6.0      |
| PoNG (ours) | <b>80.3</b> (± 4.3) | $96.9 (\pm 0.4)$  | -16.6    | PoNG (ours) | <b>79.3</b> $(\pm 0.7)$    | $98.2 (\pm 0.1)$   | -18.9      |
|             |                     |                   |          |             |                            |                    |            |

#### 



Figure 12: Transfer learning. Mean and standard deviation of test accuracy on I-RAVEN-Mesh across three random seeds. Models were trained in two setups: 1) from scratch on I-RAVEN-Mesh with variable sample size; 2) pre-trained on full I-RAVEN and fine-tuned on I-RAVEN-Mesh with variable sample size. Results for setups 1) and 2) are shown below and above the plot lines, resp. This figure extends Fig. 6 with all models considered in the paper.

Table 15: A/Size.

1116 1117 1118

1119

Table 16: A/Type.

|             | Test               | Val                | Test – Val |             | Test                | Val                 | Test – Val |
|-------------|--------------------|--------------------|------------|-------------|---------------------|---------------------|------------|
| ALANS       | $23.3 (\pm 6.5)$   | $24.6 (\pm 10.6)$  | - 1.2      | ALANS       | $19.0 (\pm 3.4)$    | $25.2 (\pm 5.2)$    | - 6.3      |
| CPCNet      | $43.5 (\pm 3.5)$   | $79.2 (\pm 2.1)$   | -35.6      | CPCNet      | $38.6 (\pm 4.3)$    | $85.2 (\pm 3.3)$    | -46.7      |
| CNN-LSTM    | $13.6 (\pm 1.4)$   | $37.2 (\pm 3.1)$   | -23.6      | CNN-LSTM    | $14.5 (\pm 0.8)$    | $38.5 (\pm 0.8)$    | -24.0      |
| CoPINet     | $21.8 (\pm 0.2)$   | $60.7 (\pm 0.1)$   | -38.9      | CoPINet     | $19.8 (\pm 0.9)$    | $58.8 (\pm 2.6)$    | -39.0      |
| DRNet       | $54.3 (\pm 3.0)$   | $93.1 (\pm 0.5)$   | -38.8      | DRNet       | $44.3 (\pm 0.8)$    | $93.9 (\pm 0.3)$    | -49.6      |
| MRNet       | $20.6 (\pm 5.0)$   | $88.1 (\pm 2.5)$   | -67.5      | MRNet       | $19.4 (\pm 0.3)$    | $93.0 (\pm 0.7)$    | -73.7      |
| PrAE        | $41.3 (\pm 1.8)$   | $59.7 (\pm 1.8)$   | -18.4      | PrAE        | $37.0 (\pm 1.7)$    | $61.1 (\pm 1.4)$    | -24.1      |
| PredRNet    | $47.5 (\pm 1.3)$   | $92.5 (\pm 0.8)$   | -45.0      | PredRNet    | $40.2 (\pm 1.3)$    | $93.9 (\pm 0.1)$    | -53.7      |
| RelBase     | $51.1 (\pm 2.4)$   | $92.9 (\pm 0.4)$   | -41.8      | RelBase     | $44.1 (\pm 1.0)$    | $93.4 (\pm 0.2)$    | -49.3      |
| SCL         | $65.6 (\pm 2.4)$   | $94.0(\pm 2.6)$    | -28.4      | SCL         | $49.5(\pm 1.8)$     | $95.9 (\pm 0.7)$    | -46.4      |
| SRAN        | $34.4 (\pm 3.0)$   | $78.1 (\pm 0.2)$   | -43.7      | SRAN        | $30.7 (\pm 2.2)$    | $78.8 (\pm 0.7)$    | -48.1      |
| STSN        | $38.4 (\pm 16.6)$  | $68.9 (\pm 34.1)$  | -30.5      | STSN        | $39.1 (\pm 5.0)$    | $66.2 \ (\pm 17.7)$ | -27.0      |
| WReN        | $12.4 \ (\pm 0.5)$ | $29.5 (\pm 0.5)$   | -17.1      | WReN        | $15.1 (\pm 0.7)$    | $23.4 (\pm 0.6)$    | -8.2       |
| PoNG (ours) | $73.5 (\pm 3.1)$   | $97.1 \ (\pm 0.5)$ | -23.6      | PoNG (ours) | <b>59.4</b> (± 6.9) | $96.1 \ (\pm 0.8)$  | -36.7      |

| Ta            | <b>ble 17:</b> A/C             | olorSize                 |               | Ta                    | ble 18: A/C                            | ColorType                 |               |
|---------------|--------------------------------|--------------------------|---------------|-----------------------|----------------------------------------|---------------------------|---------------|
|               | Test                           | Val                      | Test – Val    |                       | Test                                   | Val                       | Test – Val    |
| ALANS         | $15.1 (\pm 3.3)$               | $16.4 (\pm 4.3)$         | - 1.3         | ALANS                 | $17.7 (\pm 3.2)$                       | $20.4 (\pm 6.2)$          | - 2.7         |
| CPCNet        | $33.0(\pm 5.3)$                | $86.1(\pm 1.4)$          | -53.1         | CPCNet                | $25.0(\pm 0.9)$                        | $84.4(\pm 2.7)$           | -59.4         |
| CNN-LSTM      | $13.4(\pm 0.9)$                | $53.1(\pm 6.1)$          | -39.8         | CNN-LSTM              | $14.7(\pm 1.7)$                        | $53.3(\pm 5.9)$           | -38.7         |
| CoPINet       | $18.3 (\pm 0.3)$               | $71.7 (\pm 0.3)$         | -53.4         | CoPINet               | $17.2 (\pm 0.1)$                       | $72.8 (\pm 0.5)$          | -55.6         |
| DRNet         | $38.3 (\pm 0.5)$               | $95.9 (\pm 0.6)$         | -57.6         | DRNet                 | $29.5 (\pm 0.5)$                       | $96.0(\pm 0.4)$           | -66.4         |
| MRNet         | $18.7 (\pm 1.1)$               | $92.8~(\pm 2.1)$         | -74.0         | MRNet                 | $20.0 (\pm 2.6)$                       | $91.2 (\pm 4.2)$          | -71.2         |
| PrAE          | $30.0 (\pm 1.1)$               | $63.2~(\pm 2.3)$         | -33.1         | PrAE                  | $26.7 (\pm 0.7)$                       | $58.8 (\pm 2.6)$          | -32.1         |
| PredRNet      | $31.0 (\pm 1.6)$               | $95.4 (\pm 0.4)$         | -64.4         | PredRNet              | $28.0 (\pm 0.7)$                       | $91.6~(\pm 0.6)$          | -63.6         |
| RelBase       | $36.6 (\pm 0.8)$               | $95.2 (\pm 0.6)$         | -58.7         | RelBase               | $29.7 (\pm 0.6)$                       | $93.0 (\pm 4.0)$          | -63.3         |
| SCL           | $\frac{40.8}{22}$ (± 3.2)      | $94.2 (\pm 1.7)$         | -53.4         | SCL                   | $\frac{32.0}{22.0}$ (± 2.3)            | $96.4 (\pm 0.4)$          | -64.4         |
| SRAN          | $22.7 (\pm 1.1)$               | $84.3 (\pm 0.7)$         | -61.7         | SRAN                  | $20.9 (\pm 0.9)$                       | $84.4 (\pm 1.9)$          | -63.5         |
| STSN          | $27.3 (\pm 4.6)$               | $84.5 (\pm 12.7)$        | -57.2         | STSN                  | $21.9 (\pm 4.6)$                       | $76.2 (\pm 22.2)$         | -54.3         |
| WReN          | $13.5 (\pm 0.1)$               | $43.4 (\pm 1.4)$         | -29.9         | WReN                  | $13.8 (\pm 0.7)$                       | $41.3 (\pm 1.0)$          | -27.5         |
| PonG (ours)   | <b>44.7</b> (± 2.1)            | <b>97.3</b> (±0.5)       | -52.6         | Pong (ours)           | <b>34.3</b> (±0.8)                     | <b>90.0</b> (±0.5)        | -62.3         |
| Ta            | able 19: A/S                   | SizeType.                |               | Table 20              | ): A/Color                             | -Progres                  | sion.         |
|               | Test                           | Val                      | Test – Val    |                       | Test                                   | Val                       | Test – Val    |
| ALANS         | 157(+32)                       | $21.6 (\pm 11.7)$        | - 59          | ALANS                 | 248 (+ 18 8)                           | $26.0(\pm 20.2)$          | - 12          |
| CPCNet        | 24.1 (+12)                     | $87.0(\pm 11.7)$         | -62.9         | CPCNet                | $50.5(\pm 10.8)$                       | $75.5(\pm 0.2)$           | -250          |
| CNN-LSTM      | $13.0 (\pm 0.1)$               | $53.6 (\pm 0.2)$         | -40.6         | CNN-LSTM              | $17.2 (\pm 1.5)$                       | $29.6 (\pm 2.7)$          | -12.4         |
| CoPINet       | $19.7 (\pm 0.7)$               | $72.8 (\pm 0.3)$         | -53.1         | CoPINet               | $35.8 (\pm 0.6)$                       | $48.9(\pm 0.5)$           | -13.1         |
| DRNet         | $31.6 (\pm 1.2)$               | $95.6 (\pm 0.2)$         | -64.0         | DRNet                 | $72.8 (\pm 1.3)$                       | $96.1 (\pm 0.5)$          | -23.3         |
| MRNet         | $28.2 (\pm 0.9)$               | <b>97.7</b> $(\pm 0.2)$  | -69.5         | MRNet                 | $34.4 (\pm 3.4)$                       | $\overline{93.0}$ (± 2.7) | -58.6         |
| PrAE          | $25.6(\pm 0.8)$                | $65.0(\pm 3.7)$          | -39.4         | PrAE                  | $62.3(\pm 0.9)$                        | $65.1(\pm 3.7)$           | -2.8          |
| PredRNet      | $27.9(\pm 0.5)$                | $94.8(\pm 0.3)$          | -66.9         | PredRNet              | $62.3(\pm 2.2)$                        | $91.0(\pm 2.5)$           | -28.7         |
| RelBase       | $31.1(\pm 1.0)$                | $94.6(\pm 1.4)$          | -63.5         | RelBase               | $73.0(\pm 1.8)$                        | $95.0(\pm 0.8)$           | -22.0         |
| SCL           | <b>33.5</b> $(\pm 0.7)$        | $97.2(\pm 0.6)$          | -63.7         | SCL                   | $75.6(\pm 10.1)$                       | $88.8(\pm 7.1)$           | -13.2         |
| SRAN          | $23.3 (\pm 0.3)$               | $\overline{88.4}$ (±0.3) | -65.1         | SRAN                  | $42.1 (\pm 2.3)$                       | $64.4 (\pm 1.8)$          | -22.4         |
| STSN          | $12.3 \ (\pm 0.1)$             | $12.5~(\pm 0.5)$         | - 0.2         | STSN                  | $39.9 (\pm 14.7)$                      | $70.7~(\pm 28.9)$         | -30.8         |
| WReN          | $14.1 \ (\pm 0.2)$             | $50.9~(\pm 0.7)$         | -36.8         | WReN                  | $18.0 \ (\pm 0.4)$                     | $18.8 (\pm 0.4)$          | - 0.7         |
| PoNG (ours)   | $\underline{32.1}~(\pm2.1)$    | $96.6~(\pm 0.9)$         | -64.5         | PoNG (ours)           | $81.4 (\pm 3.1)$                       | <b>97.4</b> $(\pm 0.3)$   | -16.0         |
|               |                                |                          |               |                       |                                        |                           |               |
| <b>T</b> 11 2 | 1 . (                          |                          |               | TT 1 1 22 -           | (~ ) -                                 |                           | -1            |
| Table 2       | I: A/COLO                      | r-Aritnme                | etic.         | Table 22: F           | /Color-L                               | lstribut                  | einree.       |
|               | Test                           | Val                      | Test – Val    |                       | Test                                   | Val                       | Test – Val    |
| ALANS         | $18.3 \ (\pm 6.6)$             | $19.2 \ (\pm 6.5)$       | - 0.9         | ALANS                 | $22.4 (\pm 7.7)$                       | $21.5\;(\pm7.0)$          | + 0.9         |
| CPCNet        | $45.9 \; (\pm 2.7)$            | $74.3 (\pm 1.8)$         | -28.4         | CPCNet                | $37.8 (\pm 0.9)$                       | $74.5~(\pm 0.5)$          | -36.7         |
| CNN-LSTM      | $17.1 \ (\pm 3.7)$             | $26.0 (\pm 4.5)$         | - 8.9         | CNN-LSTM              | $20.6 (\pm 6.7)$                       | $24.9 \ (\pm 0.5)$        | - 4.3         |
| CoPINet       | $35.2 (\pm 0.5)$               | $45.5~(\pm 0.8)$         | -10.4         | CoPINet               | $26.9 (\pm 0.5)$                       | $48.7~(\pm 0.5)$          | -21.8         |
| DRNet         | $\underline{66.7}_{(\pm 1.2)}$ | $93.5 (\pm 0.5)$         | -26.8         | DRNet                 | $63.2 (\pm 0.3)$                       | $95.1 (\pm 0.2)$          | -31.9         |
| MRNet         | $35.7 (\pm 5.9)$               | $88.8 (\pm 5.6)$         | -53.1         | MRNet                 | $18.6 (\pm 0.1)$                       | $92.6~(\pm 0.4)$          | -74.0         |
| PrAE          | $43.0 (\pm 26.5)$              | $43.4(\pm 26.7)$         | - 0.4         | PrAE                  | $55.1 (\pm 0.8)$                       | $63.4 (\pm 3.7)$          | - 8.3         |
| PredRNet      | $56.9(\pm 1.4)$                | $89.5 (\pm 0.4)$         | -32.6         | PredRNet              | $48.5 (\pm 0.9)$                       | $89.7 (\pm 0.9)$          | -41.2         |
| KelBase       | $66.2 (\pm 1.0)$               | $93.5 (\pm 0.2)$         | -27.3         | RelBase               | $\frac{05.7}{62.0}$ (± 4.6)            | $93.7 (\pm 1.6)$          | -28.0         |
| SCL           | $00.0 (\pm 4.1)$               | $85.0(\pm 1.5)$          | -25.1         | SUL                   | $03.9(\pm 4.3)$                        | $83.4 (\pm 1.4)$          | -19.5         |
| SKAN<br>STEN  | $39.9(\pm 2.7)$                | $01.1 (\pm 3.2)$         | -21.8         | SKAN                  | $34.0 (\pm 3.6)$                       | $04.1 (\pm 2.3)$          | -30.1         |
| SI SIN        | $23.1 (\pm 10.6)$              | $39.0 (\pm 25.8)$        | -13.7         | SISN<br>WDoN          | $20.1 (\pm 7.7)$                       | $33.3 (\pm 23.2)$         | -32.8         |
| PoNG (ours)   | <b>700</b> $(\pm 4.1)$         | $10.2 (\pm 0.3)$         | - 1.1<br>26.2 | W KelN<br>DoNG (ours) | $11.1 (\pm 0.6)$<br><b>813</b> (± 1.0) | $10.0 (\pm 0.6)$          | - 0.0<br>15 5 |
| ong (ours)    | $/0.0 (\pm 4.1)$               | 90.3 (±0.3)              | -26.3         | POING (ours)          | <b>01.3</b> (±1.6)                     | 90.8 (±0.9)               | -15.5         |



| le 2 | I: A/Color          | r-Arithme               | tic.       | Table 22: A | /Color-L         | )ıstrıbut           | eThree.    |
|------|---------------------|-------------------------|------------|-------------|------------------|---------------------|------------|
|      | Test                | Val                     | Test – Val |             | Test             | Val                 | Test - Val |
|      | $18.3 (\pm 6.6)$    | $19.2 (\pm 6.5)$        | - 0.9      | ALANS       | $22.4 (\pm 7.7)$ | $21.5 (\pm 7.0)$    | + 0.9      |
|      | $45.9(\pm 2.7)$     | $74.3 (\pm 1.8)$        | -28.4      | CPCNet      | $37.8 (\pm 0.9)$ | $74.5 (\pm 0.5)$    | -36.7      |
| ΓМ   | $17.1 (\pm 3.7)$    | $26.0 (\pm 4.5)$        | -8.9       | CNN-LSTM    | $20.6 (\pm 6.7)$ | $24.9 \ (\pm 0.5)$  | - 4.3      |
|      | $35.2 (\pm 0.5)$    | $45.5 (\pm 0.8)$        | -10.4      | CoPINet     | $26.9 (\pm 0.5)$ | $48.7 (\pm 0.5)$    | -21.8      |
|      | $66.7(\pm 1.2)$     | $93.5(\pm 0.5)$         | -26.8      | DRNet       | $63.2 (\pm 0.3)$ | $95.1 (\pm 0.2)$    | -31.9      |
|      | $35.7 (\pm 5.9)$    | $88.8 (\pm 5.6)$        | -53.1      | MRNet       | $18.6 (\pm 0.1)$ | $92.6~(\pm 0.4)$    | -74.0      |
|      | $43.0 (\pm 26.5)$   | $43.4\;(\pm26.7)$       | -0.4       | PrAE        | $55.1 (\pm 0.8)$ | $63.4 (\pm 3.7)$    | -8.3       |
| t    | $56.9(\pm 1.4)$     | $89.5 (\pm 0.4)$        | -32.6      | PredRNet    | $48.5 (\pm 0.9)$ | $89.7~(\pm 0.9)$    | -41.2      |
|      | $66.2 (\pm 1.0)$    | $93.5 (\pm 0.2)$        | -27.3      | RelBase     | $65.7 (\pm 4.6)$ | $93.7 (\pm 1.6)$    | -28.0      |
|      | $60.0(\pm 4.1)$     | $85.0 (\pm 1.5)$        | -25.1      | SCL         | $63.9 (\pm 4.3)$ | $83.4 (\pm 1.4)$    | -19.5      |
|      | $39.9 (\pm 2.7)$    | $61.7 (\pm 3.2)$        | -21.8      | SRAN        | $34.6 (\pm 3.6)$ | $64.7 (\pm 2.3)$    | -30.1      |
|      | $25.7 (\pm 10.6)$   | $39.5~(\pm 25.8)$       | -13.7      | STSN        | $20.7 (\pm 7.7)$ | $53.5~(\pm 23.2)$   | -32.8      |
|      | 17.1 (±0.2)         | $18.2 (\pm 0.3)$        | -1.1       | WReN        | $17.7 (\pm 0.6)$ | $18.3 (\pm 0.6)$    | - 0.6      |
| ırs) | <b>70.0</b> (± 4.1) | <b>96.3</b> $(\pm 0.3)$ | -26.3      | PoNG (ours) | $81.3 (\pm 1.6)$ | <b>96.8</b> (± 0.9) | -15.5      |
|      |                     |                         |            |             |                  |                     |            |

Table 23: I-RAVEN. Results from Table 2 extended to each matrix configuration.

|             | Mean                | Center                      | 2x2Grid             | 3x3Grid            | L-R                     | U-D                 | O-IC              | O-IG                       |
|-------------|---------------------|-----------------------------|---------------------|--------------------|-------------------------|---------------------|-------------------|----------------------------|
| ALANS       | $27.0 (\pm 8.4)$    | $28.8 (\pm 9.2)$            | $25.6 (\pm 5.2)$    | $26.7 (\pm 7.7)$   | $32.1 (\pm 13.2)$       | 31.9 (± 12.0)       | $24.3 (\pm 8.5)$  | $19.8 (\pm 3.8)$           |
| CPCNet      | $70.4 (\pm 6.4)$    | $85.0 (\pm 10.9)$           | $53.1 (\pm 10.9)$   | $45.2 (\pm 7.0)$   | $89.3 (\pm 4.4)$        | $89.9(\pm 5.1)$     | $84.3 (\pm 1.3)$  | $46.0 (\pm 5.7)$           |
| CNN-LSTM    | $27.5 (\pm 1.5)$    | $41.2 (\pm 4.3)$            | $27.1 (\pm 3.1)$    | $24.8 (\pm 3.3)$   | $23.4 (\pm 1.9)$        | $22.5 (\pm 0.2)$    | $28.3 (\pm 2.0)$  | $25.6 (\pm 0.7)$           |
| CoPINet     | $43.2 (\pm 0.1)$    | $51.8 (\pm 0.7)$            | $34.0 (\pm 0.6)$    | $29.9(\pm 0.6)$    | $47.2 (\pm 0.6)$        | $49.9 (\pm 0.6)$    | $49.2 (\pm 0.9)$  | $40.3 (\pm 1.2)$           |
| DRNet       | $90.9(\pm 1.1)$     | $99.3 (\pm 0.2)$            | $91.2 (\pm 1.5)$    | $83.6 (\pm 1.1)$   | $95.8 (\pm 0.8)$        | $98.2 (\pm 0.4)$    | $98.1 (\pm 0.7)$  | $\underline{70.3}$ (± 3.8) |
| MRNet       | $86.7 (\pm 2.3)$    | <b>99.8</b> (±0.1)          | $82.2 \ (\pm 10.0)$ | $72.3\;(\pm11.1)$  | $97.8 (\pm 2.0)$        | $97.9 (\pm 1.3)$    | $94.4 (\pm 3.0)$  | $62.9 (\pm 0.8)$           |
| PrAE        | $19.5 (\pm 0.4)$    | $20.2 (\pm 0.7)$            | $36.8 (\pm 1.0)$    | $17.3 (\pm 1.8)$   | $16.4 (\pm 0.7)$        | $15.3 (\pm 0.7)$    | $15.8 (\pm 1.4)$  | $14.6 (\pm 0.6)$           |
| PredRNet    | $88.8 (\pm 1.8)$    | $98.7 (\pm 0.2)$            | $93.4(\pm 1.3)$     | $80.3 (\pm 4.8)$   | $98.5(\pm 0.5)$         | $97.7 (\pm 0.3)$    | $97.7 (\pm 1.0)$  | $55.3 (\pm 7.3)$           |
| RelBase     | $89.6 (\pm 0.6)$    | $99.1 (\pm 0.2)$            | $90.6 (\pm 0.7)$    | $82.2 (\pm 1.4)$   | $95.2 (\pm 0.3)$        | $97.7 (\pm 0.5)$    | $98.0 (\pm 0.5)$  | $64.2 (\pm 2.1)$           |
| SCL         | $83.4 (\pm 2.5)$    | $98.6 (\pm 0.2)$            | $81.6 (\pm 3.3)$    | $73.1 \ (\pm 0.7)$ | $86.7 (\pm 5.3)$        | $86.3 (\pm 5.7)$    | $88.4 (\pm 3.2)$  | $69.4 \ (\pm 0.5)$         |
| SRAN        | $58.2 (\pm 1.6)$    | $80.7 (\pm 2.8)$            | $46.7 (\pm 0.8)$    | $39.7 (\pm 1.2)$   | $68.1 (\pm 2.9)$        | $66.5 (\pm 3.1)$    | $62.6 (\pm 0.3)$  | $42.9 (\pm 1.0)$           |
| STSN        | $59.0 \ (\pm 18.5)$ | $76.1 (\pm 22.7)$           | $53.7 (\pm 15.9)$   | $48.5 (\pm 13.7)$  | $65.1 (\pm 24.5)$       | $64.1 \ (\pm 24.8)$ | $65.5~(\pm 22.3)$ | $41.7 (\pm 8.5)$           |
| WReN        | $18.4 (\pm 0.0)$    | $24.5 (\pm 3.2)$            | $17.2 (\pm 1.1)$    | $17.3 (\pm 0.9)$   | $15.1 (\pm 1.0)$        | $16.8 (\pm 1.0)$    | $18.7 (\pm 0.5)$  | $19.6 (\pm 0.4)$           |
| PoNG (ours) | $95.9 (\pm 0.7)$    | $\underline{99.5}~(\pm0.1)$ | <b>97.8</b> (± 0.9) | $91.2~(\pm 1.2)$   | <b>98.7</b> $(\pm 0.5)$ | <b>98.7</b> (± 0.6) | $98.8 (\pm 0.3)$  | $86.5 (\pm 2.1)$           |

Table 24: I-RAVEN-Mesh. Results from Table 2 extended to each matrix configuration.

| 09 |             |                     |                     |                              |                             |                         |                              | C                   |                            |
|----|-------------|---------------------|---------------------|------------------------------|-----------------------------|-------------------------|------------------------------|---------------------|----------------------------|
| 10 |             | Mean                | Center              | 2x2Grid                      | 3x3Grid                     | L-R                     | U-D                          | O-IC                | O-IG                       |
| 11 | ALANS       | $15.9 (\pm 2.6)$    | $16.3 (\pm 2.8)$    | $15.8 (\pm 2.9)$             | $16.1 (\pm 2.7)$            | $16.0 (\pm 3.5)$        | $17.1 (\pm 3.6)$             | $15.6 (\pm 1.9)$    | $14.2 (\pm 0.6)$           |
| 2  | CPCNet      | $66.6 (\pm 5.1)$    | $73.2 (\pm 11.6)$   | $53.4 (\pm 5.0)$             | $50.1 (\pm 2.9)$            | $78.4(\pm 3.9)$         | $76.9(\pm 4.4)$              | $74.9(\pm 3.2)$     | $58.9(\pm 5.4)$            |
| 2  | CNN-LSTM    | $28.9 (\pm 0.4)$    | $30.5 (\pm 0.9)$    | $27.4 \ (\pm 0.7)$           | $28.4 \ (\pm 0.5)$          | $28.9 \ (\pm 0.5)$      | $29.0\;(\pm0.8)$             | $28.8 (\pm 0.3)$    | $29.1 \ (\pm 1.5)$         |
| ,  | CoPINet     | $41.1 (\pm 0.3)$    | $41.5 (\pm 0.3)$    | $38.2 \ (\pm 0.3)$           | $34.6\;(\pm0.2)$            | $42.2 \ (\pm 1.7)$      | $42.4~(\pm 0.2)$             | $46.0 (\pm 0.3)$    | $42.9 \ (\pm 0.5)$         |
|    | DRNet       | $83.9 (\pm 2.7)$    | <b>93.7</b> (± 1.0) | $77.5 (\pm 6.6)$             | $71.4 (\pm 4.2)$            | $88.2 \ (\pm 0.8)$      | $90.9~(\pm 1.3)$             | $89.6 (\pm 0.9)$    | $\underline{76.0}$ (± 5.3) |
|    | MRNet       | $79.5 (\pm 2.0)$    | $93.2(\pm 4.9)$     | $65.4 \ (\pm 6.2)$           | $59.4~(\pm 5.1)$            | $90.1(\pm 4.2)$         | $\underline{91.2}\;(\pm3.8)$ | $89.5 (\pm 1.6)$    | $67.6 (\pm 1.7)$           |
|    | PrAE        | $33.2 (\pm 0.4)$    | $38.1 (\pm 1.0)$    | $39.1 (\pm 0.1)$             | $19.9 (\pm 0.9)$            | $41.3 (\pm 1.7)$        | $41.7 (\pm 0.7)$             | $28.4 (\pm 1.5)$    | $23.7 (\pm 1.0)$           |
|    | PredRNet    | $59.2 (\pm 6.4)$    | $67.9 (\pm 5.4)$    | $50.7 (\pm 4.3)$             | $47.2 (\pm 2.9)$            | $65.2 (\pm 8.9)$        | $63.0 (\pm 9.3)$             | $65.1 (\pm 9.5)$    | $55.0 (\pm 5.0)$           |
|    | RelBase     | $84.9(\pm 4.4)$     | $92.5 (\pm 2.3)$    | $\underline{81.5} (\pm 5.8)$ | $\underline{75.5}(\pm 5.8)$ | $88.4(\pm 3.4)$         | $90.4 \ (\pm 2.8)$           | $90.3 (\pm 5.0)$    | $75.3 (\pm 6.7)$           |
|    | SCL         | $80.9(\pm 1.5)$     | $88.5 (\pm 2.5)$    | $77.6 (\pm 1.2)$             | $73.5 (\pm 1.8)$            | $83.1 (\pm 0.0)$        | $82.8 (\pm 0.6)$             | $86.7 (\pm 0.8)$    | $74.1 (\pm 5.5)$           |
|    | SRAN        | $57.8 (\pm 0.2)$    | $65.7 (\pm 0.3)$    | $46.6 \ (\pm 0.9)$           | $45.1~(\pm 0.9)$            | $64.0~(\pm 0.2)$        | $66.1\;(\pm1.3)$             | $63.8 (\pm 0.4)$    | $53.5~(\pm 0.6)$           |
|    | STSN        | $48.7 \ (\pm 11.5)$ | $63.8~(\pm 20.2)$   | $45.2 (\pm 9.7)$             | $43.4 \ (\pm 7.3)$          | $43.9\;(\pm6.9)$        | $44.0~(\pm 7.9)$             | $52.8 \ (\pm 17.2)$ | $49.0 (\pm 12.6)$          |
|    | WReN        | $25.7 (\pm 0.2)$    | $26.4 (\pm 0.5)$    | $24.5 (\pm 0.4)$             | $25.7 \ (\pm 0.6)$          | $25.7 \ (\pm 0.3)$      | $24.9\;(\pm1.4)$             | $25.9(\pm 0.1)$     | $26.3 (\pm 0.2)$           |
|    | PoNG (ours) | 89.3 (± 2.4)        | $91.6 (\pm 3.5)$    | <b>89.6</b> (± 2.3)          | 82.8 $(\pm 2.7)$            | <b>90.8</b> $(\pm 2.4)$ | <b>91.7</b> $(\pm 2.3)$      | <b>91.0</b> (±1.8)  | 87.7 (± 2.2)               |

Table 25: A/Color. Results from Table 2 extended to each matrix configuration.

|             | Mean                | Center                     | 2x2Grid            | 3x3Grid                      | L-R                | U-D                        | O-IC                        | O-IG                        |
|-------------|---------------------|----------------------------|--------------------|------------------------------|--------------------|----------------------------|-----------------------------|-----------------------------|
| ALANS       | $ 15.2(\pm 1.4) $   | $14.7 (\pm 2.3)$           | $16.7 (\pm 2.0)$   | $15.1 (\pm 1.5)$             | $15.6 (\pm 2.8)$   | $15.7 (\pm 2.0)$           | $14.6 (\pm 1.8)$            | $14.0 (\pm 1.5)$            |
| CPCNet      | $51.2 (\pm 3.8)$    | $48.3 (\pm 6.8)$           | $41.1 (\pm 6.7)$   | $38.9 (\pm 2.6)$             | $57.4 (\pm 2.7)$   | $57.5 (\pm 2.0)$           | $67.1 (\pm 2.4)$            | $48.0 (\pm 5.7)$            |
| CNN-LSTM    | $17.0 (\pm 3.1)$    | $17.9 (\pm 4.6)$           | $17.3 (\pm 2.7)$   | $16.3 (\pm 2.3)$             | $15.8 (\pm 1.8)$   | $15.9 (\pm 2.6)$           | $18.3 (\pm 4.6)$            | $17.9 (\pm 3.6)$            |
| CoPINet     | $32.5 (\pm 0.2)$    | $33.0 \ (\pm 0.8)$         | $28.9 (\pm 0.7)$   | $27.0 \ (\pm 0.2)$           | $30.0 (\pm 0.4)$   | $30.7~(\pm 0.7)$           | $39.2 (\pm 0.6)$            | $38.7 (\pm 1.7)$            |
| DRNet       | $70.0 (\pm 1.6)$    | $66.9 (\pm 4.0)$           | $65.9 \ (\pm 1.2)$ | $\underline{63.9}\;(\pm1.1)$ | $68.1(\pm 1.7)$    | $\underline{70.1}$ (± 3.2) | $\underline{82.7}(\pm 2.0)$ | $72.5 (\pm 1.6)$            |
| MRNet       | $33.6 (\pm 8.2)$    | $27.1 (\pm 7.1)$           | $39.6 (\pm 1.4)$   | $39.4 (\pm 3.0)$             | $21.1 (\pm 16.4)$  | $20.3\;(\pm17.3)$          | $39.8 \ (\pm 16.9)$         | $48.4 (\pm 5.0)$            |
| PrAE        | $47.9 (\pm 0.9)$    | $50.0 (\pm 1.5)$           | $57.7 (\pm 1.8)$   | $36.7 (\pm 2.2)$             | $61.8 (\pm 1.6)$   | $60.8 (\pm 1.1)$           | $37.8 (\pm 0.5)$            | $30.3 (\pm 1.5)$            |
| PredRNet    | $59.4 (\pm 1.0)$    | $52.9 (\pm 1.0)$           | $61.3 \ (\pm 0.3)$ | $56.9 (\pm 2.2)$             | $58.0 \ (\pm 0.7)$ | $57.8 (\pm 1.3)$           | $69.9 (\pm 0.3)$            | $59.5 (\pm 6.0)$            |
| RelBase     | $67.4 (\pm 2.7)$    | $62.8 (\pm 4.5)$           | $66.8 (\pm 2.6)$   | $63.4 (\pm 2.4)$             | $66.4(\pm 3.4)$    | $66.2 (\pm 3.8)$           | $76.2 (\pm 3.2)$            | $70.2 (\pm 2.4)$            |
| SCL         | $65.1 (\pm 2.0)$    | $\underline{71.3}$ (± 5.2) | $66.2 (\pm 2.1)$   | $57.6 (\pm 1.9)$             | $57.5 (\pm 5.0)$   | $56.6 (\pm 4.9)$           | $74.0 (\pm 0.7)$            | $\underline{73.5}(\pm 0.6)$ |
| SRAN        | $38.3 (\pm 1.0)$    | $40.1 (\pm 2.4)$           | $35.0 \ (\pm 0.7)$ | $31.5 (\pm 1.4)$             | $35.9(\pm 3.5)$    | $36.2 (\pm 2.6)$           | $47.4 (\pm 0.7)$            | $42.0 (\pm 0.6)$            |
| STSN        | $39.3 (\pm 6.9)$    | $39.5 (\pm 1.8)$           | $39.9 (\pm 7.6)$   | $38.6 (\pm 7.0)$             | $31.9 (\pm 9.1)$   | $30.0 (\pm 7.4)$           | $49.6 (\pm 8.0)$            | $45.3 (\pm 9.0)$            |
| WReN        | $16.9 (\pm 0.5)$    | $18.7 (\pm 0.6)$           | $17.1 (\pm 0.4)$   | $16.2 (\pm 0.7)$             | $15.3 (\pm 0.5)$   | $15.7 (\pm 1.3)$           | $17.7 (\pm 0.3)$            | $17.5 (\pm 0.8)$            |
| PoNG (ours) | <b>80.3</b> (± 4.3) | $84.4 (\pm 10.1)$          | $85.4 \ (\pm 6.8)$ | $80.3 \ (\pm 5.1)$           | $72.3 (\pm 3.7)$   | $71.3 (\pm 3.8)$           | $88.9 (\pm 3.1)$            | <b>79.0</b> $(\pm 3.7)$     |
| -           |                     |                            |                    |                              |                    |                            |                             |                             |

Table 26: A/Position. Results from Table 2 extended to each matrix configuration.

| $ \begin{array}{ c c c c c c c c c c c c c c c c c c c$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |             |                   |                     |                             |                              |                   |                    |                     |                   |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------------|-------------------|---------------------|-----------------------------|------------------------------|-------------------|--------------------|---------------------|-------------------|
| $ \begin{array}{ c c c c c c c c c c c c c c c c c c c$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   |             | Mean              | Center              | 2x2Grid                     | 3x3Grid                      | L-R               | U-D                | O-IC                | O-IG              |
| $ \begin{array}{c} \mbox{CPCNet} & 68.3 (\pm 4.0) & 89.4 (\pm 8.8) & 29.1 (\pm 2.3) & 30.1 (\pm 3.6) & 96.1 (\pm 3.1) & 96.5 (\pm 2.8) & 94.3 (\pm 2.6) & 42.6 (\pm 6.6) \\ \mbox{CNN-LSTM} & 24.0 (\pm 2.9) & 43.8 (\pm 5.6) & 13.8 (\pm 1.8) & 14.5 (\pm 1.2) & 25.4 (\pm 1.9) & 26.0 (\pm 2.9) & 29.1 (\pm 5.1) & 15.1 (\pm 3.8) \\ \mbox{CoPINet} & 41.3 (\pm 1.6) & 51.2 (\pm 2.6) & 22.2 (\pm 0.7) & 22.2 (\pm 0.7) & 53.1 (\pm 2.4) & 52.8 (\pm 2.7) & 53.3 (\pm 1.6) & 34.5 (\pm 1.8) \\ \mbox{DRNet} & \frac{77.5}{1.5} (\pm 0.9) & \frac{99.4}{1.2} (\pm 0.1) & 41.8 (\pm 3.3) & 41.3 (\pm 1.6) & 97.3 (\pm 0.6) & 98.6 (\pm 0.3) & 99.2 (\pm 0.1) & 65.4 (\pm 1.1) \\ \mbox{MRNet} & 62.6 (\pm 2.6) & 92.0 (\pm 13.7) & 19.5 (\pm 3.2) & 18.2 (\pm 4.5) & 98.1 (\pm 2.5) & 97.1 (\pm 3.7) & 96.4 (\pm 4.1) & 16.9 (\pm 0.6) \\ \mbox{PrAE} & 68.2 (\pm 3.3) & 86.0 (\pm 5.5) & 56.5 (\pm 1.9) & 50.2 (\pm 1.9) & 92.0 (\pm 4.7) & 92.3 (\pm 4.5) & 62.0 (\pm 5.1) & 38.1 (\pm 1.4) \\ \mbox{PredRNet} & 73.7 (\pm 0.7) & 99.2 (\pm 0.2) & 36.2 (\pm 1.4) & 36.2 (\pm 1.5) & 98.7 (\pm 0.7) & 98.9 (\pm 0.1) & 98.6 (\pm 0.2) & 48.2 (\pm 1.4) \\ \mbox{RelBase} & 76.6 (\pm 0.3) & 99.1 (\pm 0.1) & 39.7 (\pm 1.1) & 39.9 (\pm 2.5) & 96.1 (\pm 0.3) & 98.0 (\pm 0.3) & 98.8 (\pm 0.0) & 64.3 (\pm 0.9) \\ \mbox{SCL} & 76.7 (\pm 7.1) & 99.2 (\pm 0.8) & 51.6 (\pm 7.6) & 33.4 (\pm 9.7) & 93.4 (\pm 8.9) & 93.9 (\pm 8.1) & 95.6 (\pm 5.3) & 60.0 (\pm 9.8) \\ \mbox{SCL} & 76.9 (\pm 0.7) & 80.0 (\pm 5.1) & 31.2 (\pm 0.6) & 30.3 (\pm 0.7) & 73.9 (\pm 1.2) & 74.0 (\pm 2.5) & 69.9 (\pm 0.6) & 39.0 (\pm 4.8) \\ \mbox{SCL} & 56.9 (\pm 0.7) & 80.0 (\pm 5.1) & 31.2 (\pm 0.6) & 30.3 (\pm 0.7) & 73.9 (\pm 1.2) & 74.0 (\pm 2.5) & 69.9 (\pm 0.6) & 39.0 (\pm 4.8) \\ \mbox{SCL} & 36.1 (\pm 19.9) & 54.1 (\pm 22.8) & 19.5 (\pm 6.6) & 20.5 (\pm 9.0) & 39.7 (\pm 32.8) & 41.7 (\pm 30.0) & 48.3 (\pm 30.2) & 30.1 (\pm 1.6) \\ \mbox{SCL} & 36.1 (\pm 19.9) & 54.1 (\pm 22.8) & 19.5 (\pm 6.6) & 20.5 (\pm 9.0) & 39.7 (\pm 32.8) & 41.7 (\pm 30.0) & 48.3 (\pm 30.2) & 30.1 (\pm 10.6) \\ \mbox{SCL} & 36.1 (\pm 19.9) & 54.1 (\pm 22.8) & 19.5 (\pm 6.6) & 20.5 (\pm 9.0) & 39.7 (\pm 32.8) & 41.7 (\pm 30.0) & 48.3 (\pm 30.2) & 30.1 (\pm 10.6) & 30.0 (\pm 5.6) & 30.0 (\pm 5.6$                                                                                                                                                           | ALANS       | $16.0 (\pm 1.0)$  | $15.0 (\pm 1.9)$    | $18.1 (\pm 0.8)$            | $16.8 (\pm 1.3)$             | $16.7 (\pm 1.6)$  | $16.5 (\pm 1.0)$   | $15.2 \ (\pm 0.5)$  | $14.0 (\pm 2.6)$  |
| $ \begin{array}{c} \text{CNN-LSTM} \\ \text{CNN-LSTM} \\ \text{24.0} (\pm 2.9) \\ 43.8 (\pm 5.6) \\ 13.8 (\pm 1.8) \\ 14.5 (\pm 1.2) \\ 25.4 (\pm 1.9) \\ 26.0 (\pm 2.9) \\ 29.1 (\pm 5.1) \\ 15.1 (\pm 3.1) \\ 15.1 (\pm 3.1) \\ 20.1 (\pm 2.6) \\ 22.2 (\pm 0.7) \\ 22.2 (\pm 0.7) \\ 22.2 (\pm 0.7) \\ 53.1 (\pm 2.4) \\ 52.8 (\pm 2.7) \\ 53.3 (\pm 1.6) \\ 34.5 (\pm 1.6)$ | CPCNet      | $68.3 (\pm 4.0)$  | $89.4 (\pm 8.8)$    | $29.1 (\pm 2.3)$            | $30.1 (\pm 3.6)$             | $96.1 (\pm 3.1)$  | $96.5 (\pm 2.8)$   | $94.3 (\pm 2.6)$    | $42.6 (\pm 6.2)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | CNN-LSTM    | $24.0 (\pm 2.9)$  | $43.8 (\pm 5.6)$    | $13.8 (\pm 1.8)$            | $14.5 (\pm 1.2)$             | $25.4 (\pm 1.9)$  | $26.0 (\pm 2.9)$   | $29.1 (\pm 5.1)$    | $15.1 (\pm 3.2)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | CoPINet     | $41.3 (\pm 1.6)$  | $51.2 (\pm 2.6)$    | $22.2 (\pm 0.7)$            | $22.2 (\pm 0.7)$             | $53.1 (\pm 2.4)$  | $52.8 (\pm 2.7)$   | $53.3 (\pm 1.6)$    | $34.5 (\pm 1.6)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | DRNet       | $77.5(\pm 0.9)$   | $99.4 (\pm 0.1)$    | $41.8 (\pm 3.3)$            | $41.3 (\pm 1.6)$             | $97.3 (\pm 0.6)$  | $98.6 (\pm 0.3)$   | <b>99.2</b> (± 0.1) | 65.4 (±1.2)       |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | MRNet       | $62.6 (\pm 2.6)$  | $92.0 (\pm 13.7)$   | $19.5 (\pm 3.2)$            | $18.2 (\pm 4.5)$             | $98.1 (\pm 2.5)$  | $97.1 (\pm 3.7)$   | $96.4(\pm 4.1)$     | $16.9 (\pm 0.8)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | PrAE        | $68.2 (\pm 3.3)$  | $86.0 (\pm 5.5)$    | 56.5 (±1.9)                 | <b>50.2</b> (±1.9)           | $92.0 (\pm 4.7)$  | $92.3 (\pm 4.5)$   | $62.0 (\pm 5.1)$    | $38.1 (\pm 1.6)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | PredRNet    | $73.7 (\pm 0.7)$  | $99.2 (\pm 0.2)$    | $36.2 (\pm 1.4)$            | $36.2 (\pm 1.5)$             | $98.7 (\pm 0.7)$  | <b>98.9</b> (±0.1) | $98.6 (\pm 0.2)$    | $48.2 (\pm 1.1)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | RelBase     | $76.6 (\pm 0.3)$  | $99.1 (\pm 0.1)$    | $39.7 (\pm 1.1)$            | $39.9 (\pm 2.5)$             | $96.1 (\pm 0.3)$  | $98.0 (\pm 0.3)$   | $98.8 (\pm 0.0)$    | $64.3 (\pm 0.8)$  |
| $ \begin{array}{llllllllllllllllllllllllllllllllllll$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     | SCL         | $76.7 (\pm 7.1)$  | $99.2 (\pm 0.8)$    | $51.6 (\pm 7.6)$            | $43.4 (\pm 9.7)$             | $93.4 (\pm 8.9)$  | $93.9(\pm 8.1)$    | $95.6 (\pm 5.3)$    | $60.0 (\pm 9.5)$  |
| $\textbf{STSN} \qquad \qquad 36.1 \ (\pm \ 19.9) \ 54.1 \ (\pm \ 22.8) \ \ 19.5 \ (\pm \ 6.6) \ \ 20.5 \ (\pm \ 9.0) \ \ 39.7 \ (\pm \ 32.8) \ \ 41.7 \ (\pm \ 30.0) \ \ 48.3 \ (\pm \ 30.2) \ \ 30.1 \ (\pm \ 10.6) \ \ 30.1 \ \ 30.1 \ \ 30.2 \ \ 30.1 \ \ 30.1 \ \ 30.2 \ \ 30.1 \ \ 30.2 \ \ 30.1 \ \ 30.2 \ \ 30.1 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ 30.2 \ \ \ 30.2 \ \ \ 30.2 \ \ \ 30.2 \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \ \$                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  | SRAN        | $56.9 (\pm 0.7)$  | $80.0 (\pm 5.1)$    | $31.2 (\pm 0.6)$            | $30.3 (\pm 0.7)$             | $73.9(\pm 1.2)$   | $74.0 (\pm 2.5)$   | $69.9 (\pm 0.6)$    | $39.0 (\pm 4.4)$  |
|                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           | STSN        | $36.1 (\pm 19.9)$ | $54.1 (\pm 22.8)$   | $19.5 (\pm 6.6)$            | $20.5 (\pm 9.0)$             | $39.7 (\pm 32.8)$ | $41.7 (\pm 30.0)$  | $48.3 (\pm 30.2)$   | $30.1 (\pm 10.6)$ |
| WReN $17.3 (\pm 0.4) 21.1 (\pm 1.5) 15.6 (\pm 0.3) 15.9 (\pm 1.2) 15.4 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 1.2) 16.8 (\pm 1.1) 16.3 (\pm 0.5) 19.7 (\pm 0.5$                                                                                                                  | WReN        | $17.3 (\pm 0.4)$  | $21.1 (\pm 1.5)$    | $15.6 (\pm 0.3)$            | $15.9(\pm 1.2)$              | $15.4 (\pm 1.1)$  | $16.3 (\pm 0.5)$   | $19.7 (\pm 1.2)$    | $16.8 (\pm 1.1)$  |
| $\begin{array}{c c c c c c c c c c c c c c c c c c c $                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    | PoNG (ours) | $79.3 (\pm 0.7)$  | <b>99.6</b> (± 0.1) | $\underline{54.5}~(\pm1.0)$ | $\underline{44.0}\;(\pm2.8)$ | $98.8 (\pm 0.1)$  | $98.9 \ (\pm 0.2)$ | $98.7 \ (\pm 0.2)$  | $60.6~(\pm 2.8)$  |

Table 27: A/Size. Results from Table 2 extended to each matrix configuration.

|             | Mean                | Center              | 2x2Grid                    | 3x3Grid             | L-R               | U-D                 | O-IC              | 0-1           |
|-------------|---------------------|---------------------|----------------------------|---------------------|-------------------|---------------------|-------------------|---------------|
| ALANS       | $23.3 (\pm 6.5)$    | $23.9(\pm 7.4)$     | $23.6 (\pm 4.9)$           | $22.1 (\pm 5.8)$    | $26.6 (\pm 8.1)$  | $26.9(\pm 8.6)$     | $22.0(\pm 7.1)$   | 18.4 (:       |
| CPCNet      | $43.5(\pm 3.5)$     | $48.8 (\pm 6.3)$    | $38.2 (\pm 3.9)$           | $36.3 (\pm 1.8)$    | $51.8 (\pm 4.1)$  | $50.6 (\pm 1.7)$    | $47.8 (\pm 6.2)$  | 31.0 (        |
| CNN-LSTM    | $13.6 (\pm 1.4)$    | $15.8 (\pm 2.4)$    | $16.6 (\pm 1.3)$           | $15.0 (\pm 2.3)$    | $12.7 (\pm 1.3)$  | $12.3 (\pm 1.1)$    | $11.2 (\pm 1.9)$  | 11.5 (        |
| CoPINet     | $21.8 (\pm 0.2)$    | $22.9(\pm 1.0)$     | $25.9 (\pm 1.2)$           | $27.0 (\pm 0.3)$    | $20.1 (\pm 0.9)$  | $22.1 (\pm 0.5)$    | $15.0 (\pm 0.5)$  | 19.4 (        |
| DRNet       | $54.3 (\pm 3.0)$    | $63.0 (\pm 4.8)$    | $56.5 (\pm 1.3)$           | $52.7 (\pm 1.5)$    | $52.4 (\pm 2.8)$  | $56.6 (\pm 3.9)$    | $48.9(\pm 7.5)$   | 50.2 (        |
| MRNet       | $20.6 (\pm 5.0)$    | $20.8 \ (\pm 12.1)$ | $37.6 (\pm 6.3)$           | $35.1 (\pm 5.7)$    | $8.1(\pm 4.1)$    | $8.7 (\pm 2.7)$     | $12.1 (\pm 5.1)$  | 21.7 (        |
| PrAE        | $41.3 (\pm 1.8)$    | $38.9(\pm 4.1)$     | $49.5 (\pm 1.6)$           | $34.4 (\pm 3.8)$    | $54.2 (\pm 1.7)$  | $53.9(\pm 3.5)$     | $31.4 (\pm 2.2)$  | 27.0 (        |
| PredRNet    | $47.5 (\pm 1.3)$    | $55.4 (\pm 3.9)$    | $52.9 (\pm 0.5)$           | $50.3 (\pm 0.3)$    | $48.2 (\pm 1.3)$  | $48.1 (\pm 1.0)$    | $39.0 (\pm 0.3)$  | 38.1 (        |
| RelBase     | $51.1 (\pm 2.4)$    | $59.3 (\pm 1.8)$    | $54.3 (\pm 0.7)$           | $52.8 (\pm 0.5)$    | $50.3 (\pm 3.0)$  | $52.8 (\pm 1.5)$    | $42.0 (\pm 7.1)$  | 46.5 (        |
| SCL         | $65.6 (\pm 2.4)$    | $66.2 (\pm 3.9)$    | $\underline{72.1}$ (± 5.5) | $67.0 (\pm 4.5)$    | $61.7 (\pm 0.5)$  | $62.0 (\pm 0.3)$    | $68.6 (\pm 3.7)$  | <u>61.5</u> ( |
| SRAN        | $34.4 (\pm 3.0)$    | $39.4 (\pm 7.7)$    | $35.9 (\pm 2.0)$           | $35.9 (\pm 1.3)$    | $35.8 (\pm 2.9)$  | $37.3 (\pm 4.6)$    | $30.4 (\pm 6.1)$  | 26.2(         |
| STSN        | $38.4 (\pm 16.6)$   | $40.6 (\pm 14.9)$   | $42.0\;(\pm21.7)$          | $38.4 (\pm 18.5)$   | $36.8 (\pm 16.8)$ | $38.3 \ (\pm 16.1)$ | $40.1~(\pm 16.7)$ | 33.1 (:       |
| WReN        | $12.4 (\pm 0.5)$    | $13.8 (\pm 0.5)$    | $14.1 (\pm 0.1)$           | $15.0 (\pm 0.6)$    | $13.3 (\pm 0.9)$  | $13.3 (\pm 0.9)$    | $8.2(\pm 1.0)$    | 9.2 (=        |
| PoNG (ours) | <b>73.5</b> (± 3.1) | 84.0 (± 2.9)        | $81.1 (\pm 8.1)$           | <b>75.2</b> (± 8.6) | 67.7 (± 3.4)      | 65.3 (± 1.1)        | 78.9 (± 3.3)      | 62.5 (        |

Table 28: A/Type. Results from Table 2 extended to each matrix configuration.

|             | Mean                | Center             | 2x2Grid             | 3x3Grid             | L-R                 | U-D                 | O-IC                    | O-IG             |
|-------------|---------------------|--------------------|---------------------|---------------------|---------------------|---------------------|-------------------------|------------------|
| ALANS       | $19.0 (\pm 3.4)$    | $18.5 (\pm 4.9)$   | $19.7 (\pm 1.7)$    | $18.9 (\pm 3.8)$    | $22.1 (\pm 3.5)$    | $20.8 (\pm 4.5)$    | $17.6 (\pm 3.1)$        | $15.2 (\pm 3.1)$ |
| CPCNet      | $38.6 (\pm 4.3)$    | $33.3 \ (\pm 4.5)$ | $45.5 (\pm 5.0)$    | $42.4~(\pm 3.1)$    | $40.8 (\pm 3.4)$    | $40.4~(\pm 4.3)$    | $32.2~(\pm 5.0)$        | $34.7 (\pm 6.6)$ |
| CNN-LSTM    | $14.5 (\pm 0.8)$    | $13.1 \ (\pm 0.7)$ | $17.0 \ (\pm 0.7)$  | $17.3 \ (\pm 0.8)$  | $13.8 (\pm 1.9)$    | $13.5 \ (\pm 0.8)$  | $13.9 \ (\pm 1.8)$      | $12.9 (\pm 1.1)$ |
| CoPINet     | $19.8 (\pm 0.9)$    | $18.4 (\pm 1.8)$   | $25.8 (\pm 0.5)$    | $25.5 \ (\pm 0.6)$  | $18.7 (\pm 1.7)$    | $16.9\;(\pm1.6)$    | $13.8 (\pm 1.0)$        | $19.5 (\pm 0.8)$ |
| DRNet       | $44.3 (\pm 0.8)$    | $38.7 (\pm 1.8)$   | $52.0 (\pm 1.5)$    | $50.2 (\pm 2.4)$    | $41.6 (\pm 0.4)$    | $41.4 (\pm 0.8)$    | $38.7 (\pm 1.7)$        | $47.3 (\pm 0.3)$ |
| MRNet       | $19.4 (\pm 0.3)$    | $2.7~(\pm 0.2)$    | $48.0 (\pm 1.8)$    | $44.9 (\pm 0.4)$    | $9.0(\pm 3.8)$      | $7.0~(\pm 0.9)$     | $6.6 (\pm 1.3)$         | $17.8 (\pm 6.0)$ |
| PrAE        | $37.0(\pm 1.7)$     | $37.1 (\pm 1.8)$   | $46.7 (\pm 1.6)$    | $30.9 (\pm 1.8)$    | $46.6 (\pm 2.1)$    | $46.9 (\pm 2.5)$    | $25.6 (\pm 4.7)$        | $25.4 (\pm 0.6)$ |
| PredRNet    | $40.2 (\pm 1.3)$    | $30.3 \ (\pm 2.3)$ | $49.1 (\pm 1.2)$    | $51.8 (\pm 2.4)$    | $39.6 (\pm 2.5)$    | $37.5~(\pm 2.2)$    | $32.8 (\pm 2.2)$        | $39.9 (\pm 0.8)$ |
| RelBase     | $44.1 (\pm 1.0)$    | $39.4 \ (\pm 0.9)$ | $50.9 (\pm 0.5)$    | $51.0\ (\pm0.4)$    | $40.6~(\pm 2.4)$    | $41.6\;(\pm0.8)$    | $37.6 \ (\pm 1.5)$      | $47.1 (\pm 1.4)$ |
| SCL         | $49.5(\pm 1.8)$     | $47.5 (\pm 2.4)$   | $58.8 (\pm 4.0)$    | $57.6 (\pm 0.6)$    | $47.8(\pm 1.1)$     | $47.2(\pm 1.6)$     | $40.9 (\pm 4.0)$        | $46.3 (\pm 2.0)$ |
| SRAN        | $30.7 (\pm 2.2)$    | $30.2~(\pm 1.0)$   | $37.3 (\pm 1.3)$    | $36.3 \ (\pm 0.9)$  | $33.0 \ (\pm 3.3)$  | $32.3~(\pm 3.5)$    | $21.3 (\pm 4.6)$        | $23.9 (\pm 1.9)$ |
| STSN        | $39.1 (\pm 5.0)$    | $48.5(\pm 3.5)$    | $44.6 \ (\pm 11.3)$ | $40.6 \ (\pm 9.7)$  | $38.3 (\pm 8.5)$    | $36.8 \ (\pm 7.7)$  | $35.1 \ (\pm 1.6)$      | $30.1 (\pm 0.9)$ |
| WReN        | $15.1 (\pm 0.7)$    | $15.5 (\pm 2.0)$   | $15.6 (\pm 0.5)$    | $16.3 (\pm 1.9)$    | $14.1 (\pm 0.1)$    | $14.8 (\pm 0.3)$    | $14.5 (\pm 1.4)$        | $15.0 (\pm 0.5)$ |
| PoNG (ours) | <b>59.4</b> (± 6.9) | 58.4 (± 7.4)       | 69.7 (± 7.9)        | <b>65.4</b> (± 7.4) | <b>56.6</b> (± 8.1) | <b>56.6</b> (± 8.9) | <b>53.2</b> $(\pm 6.3)$ | 55.4 (± 4.2      |

 Table 29: A/ColorSize. The table presents results from Table 10 extended to each matrix configuration.

|             | Mean                | Center             | 2x2-Grid            | 3x3-Grid           | L-R                | U-D                | O-IC                    | O-IG           |
|-------------|---------------------|--------------------|---------------------|--------------------|--------------------|--------------------|-------------------------|----------------|
| ALANS       | $ 15.1(\pm 3.3) $   | $15.2 (\pm 4.3)$   | $17.2 (\pm 2.7)$    | $17.2 (\pm 3.9)$   | $14.9 (\pm 3.4)$   | $13.8 (\pm 4.5)$   | $13.7 (\pm 1.9)$        | 13.7 (± 2      |
| CPCNet      | $33.0(\pm 5.3)$     | $33.0 (\pm 4.4)$   | $34.6 (\pm 4.7)$    | $32.6 (\pm 2.2)$   | $28.0 (\pm 5.7)$   | $30.8 \ (\pm 6.5)$ | $35.8 (\pm 7.6)$        | $36.3 (\pm$    |
| CNN-LSTM    | $13.4 (\pm 0.9)$    | $14.3 (\pm 1.0)$   | $15.6 (\pm 2.4)$    | $14.4 (\pm 1.0)$   | $13.4 \ (\pm 0.3)$ | $13.8 (\pm 1.1)$   | $11.7 (\pm 1.7)$        | $10.4 (\pm$    |
| CoPINet     | $18.3 (\pm 0.3)$    | $19.7 \ (\pm 0.8)$ | $24.8 (\pm 0.3)$    | $23.5 \ (\pm 2.0)$ | $13.9 \ (\pm 1.3)$ | $15.2 \ (\pm 0.3)$ | $12.5~(\pm 0.5)$        | $18.4 (\pm$    |
| DRNet       | $38.3 (\pm 0.5)$    | $38.3 (\pm 0.9)$   | $44.2 (\pm 2.9)$    | $39.5 (\pm 2.2)$   | $33.0(\pm 1.1)$    | $35.2 (\pm 3.3)$   | $37.2 (\pm 3.6)$        | $40.2(\pm$     |
| MRNet       | $18.7 (\pm 1.1)$    | $17.6 \ (\pm 0.9)$ | $26.5 (\pm 2.6)$    | $26.7 (\pm 1.8)$   | $13.8 (\pm 1.6)$   | $13.9 \ (\pm 1.1)$ | $14.8 (\pm 1.6)$        | $18.1 (\pm$    |
| PrAE        | $30.0(\pm 1.1)$     | $27.1 (\pm 2.6)$   | $40.1 (\pm 1.2)$    | $30.8 (\pm 1.9)$   | $29.9(\pm 1.3)$    | $30.6 (\pm 1.2)$   | $26.2 (\pm 1.5)$        | $25.7 (\pm$    |
| PredRNet    | $31.0(\pm 1.6)$     | $31.2 (\pm 2.4)$   | $40.5 (\pm 0.8)$    | $37.0 (\pm 2.3)$   | $25.9 (\pm 1.6)$   | $27.9 (\pm 2.9)$   | $23.7 (\pm 0.5)$        | $30.0(\pm$     |
| RelBase     | $36.6 (\pm 0.8)$    | $36.5 (\pm 0.5)$   | $41.4 (\pm 0.3)$    | $38.4 (\pm 1.3)$   | $32.2 (\pm 0.4)$   | $37.1 (\pm 0.6)$   | $33.2 (\pm 3.2)$        | $36.8(\pm$     |
| SCL         | $40.8(\pm 3.2)$     | $40.8 (\pm 2.4)$   | $49.2(\pm 8.2)$     | $45.7 (\pm 5.5)$   | $29.8 (\pm 3.3)$   | $30.7 (\pm 2.0)$   | $43.4(\pm 2.9)$         | 46.0 (±        |
| SRAN        | $22.7 (\pm 1.1)$    | $18.7 (\pm 2.1)$   | $29.3 (\pm 1.8)$    | $26.6 (\pm 0.5)$   | $19.0 (\pm 2.4)$   | $18.4 (\pm 2.5)$   | $22.5 (\pm 1.9)$        | $23.7(\pm$     |
| STSN        | $27.3 (\pm 4.6)$    | $28.9 (\pm 6.1)$   | $30.3 (\pm 5.3)$    | $29.3 (\pm 5.9)$   | $24.3 (\pm 6.0)$   | $24.6 (\pm 5.0)$   | $28.8 (\pm 5.8)$        | 26.1 (±        |
| WReN        | $13.5 (\pm 0.1)$    | $14.6 (\pm 0.7)$   | $14.3 (\pm 0.2)$    | $15.0 (\pm 0.4)$   | $15.2 (\pm 0.2)$   | $14.4 (\pm 0.8)$   | $11.1 (\pm 1.2)$        | $10.1 (\pm$    |
| PoNG (ours) | <b>44.7</b> (± 2.1) | $46.1 (\pm 3.1)$   | <b>53.5</b> (± 2.5) | $48.4 (\pm 3.4)$   | $36.9 (\pm 3.1)$   | $35.1 \ (\pm 2.9)$ | <b>48.6</b> $(\pm 4.0)$ | <u>44.2</u> (± |

Table 30: A/ColorType. The table presents results from Table 10 extended to each matrix con-<br/>figuration.

|             | Mean             | Center           | 2x2-Grid                     | 3x3-Grid           | L-R                         | U-D                | O-IC               | O-IG             |
|-------------|------------------|------------------|------------------------------|--------------------|-----------------------------|--------------------|--------------------|------------------|
| ALANS       | $17.7 (\pm 3.2)$ | $17.4(\pm 4.4)$  | $18.4 (\pm 1.6)$             | $17.6 (\pm 2.4)$   | $19.7 (\pm 5.9)$            | $19.0 (\pm 5.4)$   | $17.1 (\pm 3.4)$   | $14.9 (\pm 1.3)$ |
| CPCNet      | $25.0(\pm 0.9)$  | $21.0(\pm 1.2)$  | $33.0(\pm 4.5)$              | $28.4(\pm 1.6)$    | $23.3(\pm 1.1)$             | $20.9(\pm 1.7)$    | $21.5(\pm 2.4)$    | $26.5(\pm 2.0)$  |
| CNN-LSTM    | $14.7 (\pm 1.7)$ | $13.9 (\pm 3.2)$ | $18.6 (\pm 1.6)$             | $16.7 (\pm 1.9)$   | $14.7 (\pm 1.3)$            | $12.8 (\pm 0.8)$   | $13.4 (\pm 3.0)$   | $12.4 (\pm 2.8)$ |
| CoPINet     | $17.2 (\pm 0.1)$ | $14.4 (\pm 0.5)$ | $23.2 (\pm 0.8)$             | $23.6 (\pm 0.7)$   | $13.6 (\pm 0.7)$            | $13.5 (\pm 0.6)$   | $13.2 (\pm 1.0)$   | $18.6 (\pm 1.0)$ |
| DRNet       | $29.5\;(\pm0.5)$ | $22.7~(\pm 2.0)$ | $36.7 (\pm 0.3)$             | $34.7 \ (\pm 0.5)$ | $24.3 (\pm 1.2)$            | $23.0 \ (\pm 1.6)$ | $27.2 \ (\pm 0.6)$ | $38.0 (\pm 0.2)$ |
| MRNet       | $20.0 (\pm 2.6)$ | $15.5 (\pm 0.6)$ | $31.3 (\pm 2.7)$             | $28.6 (\pm 0.7)$   | $15.0 (\pm 3.9)$            | $14.6 (\pm 4.4)$   | $15.1 (\pm 3.6)$   | $19.7 (\pm 3.7)$ |
| PrAE        | $26.7 (\pm 0.7)$ | $24.3 (\pm 1.0)$ | $35.9(\pm 0.5)$              | $27.3 (\pm 1.6)$   | <b>30.3</b> (± 1.2)         | $28.1 (\pm 0.9)$   | $17.5 (\pm 0.6)$   | $23.1 (\pm 1.1)$ |
| PredRNet    | $28.0 (\pm 0.7)$ | $20.6 (\pm 1.1)$ | $38.3(\pm 1.6)$              | $35.2 (\pm 1.6)$   | $26.1 (\pm 1.4)$            | $23.8 (\pm 0.6)$   | $23.3(\pm 1.1)$    | $28.7 (\pm 1.5)$ |
| RelBase     | $29.7 (\pm 0.6)$ | $23.2 (\pm 0.8)$ | $36.6 (\pm 0.6)$             | $34.4 (\pm 1.0)$   | $25.0 (\pm 2.0)$            | $23.7 (\pm 0.9)$   | $27.5(\pm 1.1)$    | $37.0 (\pm 0.2)$ |
| SCL         | $32.0(\pm 2.3)$  | $25.9 (\pm 3.8)$ | <b>44.4</b> $(\pm 2.4)$      | $40.1(\pm 1.9)$    | $26.1 (\pm 3.0)$            | $24.4 (\pm 1.8)$   | $27.1 (\pm 1.9)$   | $35.7 (\pm 2.1)$ |
| SRAN        | $20.9 (\pm 0.9)$ | $20.2 (\pm 2.6)$ | $27.2 (\pm 1.9)$             | $24.9 (\pm 1.9)$   | $19.2 (\pm 0.6)$            | $18.2 (\pm 0.9)$   | $15.7 (\pm 0.5)$   | $20.5 (\pm 0.7)$ |
| STSN        | $21.9(\pm 4.6)$  | $21.9(\pm 4.1)$  | $26.9(\pm 7.7)$              | $25.3(\pm 8.1)$    | $20.5 (\pm 2.0)$            | $20.7 (\pm 2.7)$   | $18.6 (\pm 4.7)$   | $19.8 (\pm 4.7)$ |
| WReN        | $13.8 (\pm 0.7)$ | $13.1 (\pm 2.0)$ | $16.0 (\pm 0.5)$             | $14.5 (\pm 0.7)$   | $13.3 (\pm 1.1)$            | $14.0(\pm 1.1)$    | $13.3 (\pm 0.3)$   | $13.1 (\pm 1.1)$ |
| PoNG (ours) | $34.3 (\pm 0.8)$ | $29.4 (\pm 1.4)$ | $\underline{43.4}\;(\pm1.4)$ | $41.2 (\pm 2.8)$   | $\underline{29.5}~(\pm1.2)$ | $29.2 (\pm 0.8)$   | $29.7 (\pm 1.6)$   | $37.2 (\pm 1.6)$ |

1334Table 31: A/SizeType. The table presents results from Table 10 extended to each matrix config-1335uration.

|             | Mean               | Center             | 2x2-Grid                | 3x3-Grid            | L-R                     | U-D                     | O-IC                    | 0-I     |
|-------------|--------------------|--------------------|-------------------------|---------------------|-------------------------|-------------------------|-------------------------|---------|
| ALANS       | 15.7 (± 3.2)       | $14.0 (\pm 2.0)$   | $18.6 (\pm 5.4)$        | $18.2 (\pm 5.7)$    | $15.1 (\pm 2.0)$        | $14.8 (\pm 1.8)$        | $14.0 (\pm 1.3)$        | 15.6 (± |
| CPCNet      | $24.1 (\pm 1.2)$   | $26.3 (\pm 1.8)$   | $30.1 (\pm 2.0)$        | $29.0 (\pm 0.4)$    | $22.3 (\pm 0.6)$        | $22.8 (\pm 0.5)$        | $17.8(\pm 1.4)$         | 20.1 (± |
| CNN-LSTM    | $13.0 (\pm 0.1)$   | $13.1 (\pm 0.4)$   | $14.2 (\pm 0.9)$        | $14.0 (\pm 0.6)$    | $13.2 (\pm 0.9)$        | $12.1 (\pm 0.0)$        | $12.3 (\pm 0.7)$        | 12.6 (± |
| CoPINet     | $19.7 (\pm 0.7)$   | $20.9 (\pm 1.6)$   | $23.5 (\pm 0.8)$        | $22.6 (\pm 0.9)$    | $19.1 (\pm 1.1)$        | $19.0 (\pm 0.7)$        | $17.4 (\pm 1.8)$        | 15.5 (± |
| DRNet       | $31.6 (\pm 1.2)$   | $35.6 (\pm 2.6)$   | $40.4 (\pm 0.7)$        | $37.1 (\pm 1.1)$    | $28.4 (\pm 3.1)$        | $27.5 (\pm 1.7)$        | $22.9 (\pm 0.7)$        | 29.4 (± |
| MRNet       | $28.2 (\pm 0.9)$   | $31.7 (\pm 2.0)$   | $34.9(\pm 2.4)$         | $31.0 (\pm 1.5)$    | $27.2 (\pm 3.5)$        | $27.1 (\pm 3.3)$        | $21.1 (\pm 1.3)$        | 24.4 (= |
| PrAE        | $25.6 \ (\pm 0.8)$ | $23.1 \ (\pm 1.1)$ | $35.2 (\pm 1.0)$        | $27.1 (\pm 2.0)$    | $28.0 \ (\pm 0.6)$      | $26.3 \ (\pm 0.9)$      | $17.6 \ (\pm 0.3)$      | 21.9 (= |
| PredRNet    | $27.9 (\pm 0.5)$   | $25.8 (\pm 1.3)$   | $38.4 (\pm 1.3)$        | $35.2 (\pm 0.6)$    | $25.0 (\pm 0.7)$        | $26.7 (\pm 0.7)$        | $17.8 (\pm 1.5)$        | 25.9 (= |
| RelBase     | $31.1 (\pm 1.0)$   | $33.3 (\pm 3.0)$   | $39.3 (\pm 0.8)$        | $37.9 (\pm 0.8)$    | $28.4 (\pm 1.9)$        | $29.2(\pm 1.9)$         | $22.1 (\pm 1.6)$        | 27.5 (= |
| SCL         | $33.5 (\pm 0.7)$   | $41.3 (\pm 0.6)$   | <b>44.0</b> $(\pm 2.4)$ | $39.7 (\pm 0.8)$    | <b>30.7</b> $(\pm 0.4)$ | <b>29.3</b> $(\pm 1.6)$ | <b>23.2</b> $(\pm 0.6)$ | 26.1 (= |
| SRAN        | $23.3\;(\pm0.3)$   | $23.7~(\pm 2.6)$   | $30.8 (\pm 1.7)$        | $28.2 \ (\pm 0.6)$  | $21.0~(\pm 2.6)$        | $22.7~(\pm 0.8)$        | $16.0 \ (\pm 0.6)$      | 20.6 (= |
| STSN        | $12.3\;(\pm0.1)$   | $12.9 \ (\pm 1.0)$ | $12.0 \ (\pm 0.8)$      | $12.4 (\pm 1.8)$    | $11.9 \ (\pm 0.2)$      | $11.9 \ (\pm 0.7)$      | $12.2 (\pm 0.6)$        | 12.3 (= |
| WReN        | $14.1 (\pm 0.2)$   | $14.3 (\pm 0.5)$   | $15.2 (\pm 0.8)$        | $14.7 (\pm 0.4)$    | $14.0 (\pm 0.8)$        | $14.7 (\pm 0.5)$        | $12.1 (\pm 0.4)$        | 13.6(:  |
| PoNG (ours) | $32.1(\pm 2.1)$    | $34.6 (\pm 3.3)$   | $42.4 (\pm 0.7)$        | <b>39.9</b> (± 2.8) | $28.8(\pm 1.8)$         | $28.4 (\pm 1.9)$        | $22.1 (\pm 4.4)$        | 28.6 (: |

| -4 | 0  | - | -4 |
|----|----|---|----|
|    | -5 |   |    |
|    | 0  | 5 |    |

Table 32: A/Color-Progression. The table presents results from Table 10 extended to each matrix configuration. 

|             | Mean              | Center            | 2x2-Grid                   | 3x3-Grid            | L-R                        | U-D                         | O-IC                | O-IG              |
|-------------|-------------------|-------------------|----------------------------|---------------------|----------------------------|-----------------------------|---------------------|-------------------|
| ALANS       | $24.8 (\pm 18.8)$ | $24.9 (\pm 21.3)$ | $24.5 (\pm 13.8)$          | $24.3 (\pm 16.4)$   | $26.3 (\pm 23.3)$          | $29.3 (\pm 26.3)$           | $24.5 (\pm 20.1)$   | $19.7 (\pm 10.1)$ |
| CPCNet      | $50.5 (\pm 0.6)$  | $51.9(\pm 0.3)$   | $37.6 (\pm 2.6)$           | $33.9 (\pm 1.5)$    | $59.0 (\pm 0.6)$           | $58.9(\pm 1.5)$             | $67.7 (\pm 2.4)$    | $44.8 (\pm 0.8)$  |
| CNN-LSTM    | $17.2 (\pm 1.5)$  | $20.3 (\pm 2.2)$  | $17.7 (\pm 1.4)$           | $16.9 (\pm 2.4)$    | $16.0 (\pm 1.3)$           | $15.6 (\pm 1.4)$            | $18.0 (\pm 1.0)$    | $15.7 (\pm 1.9)$  |
| CoPINet     | $35.8 (\pm 0.6)$  | $37.3 (\pm 0.9)$  | $29.8 (\pm 1.8)$           | $28.0 (\pm 0.3)$    | $35.8 (\pm 0.5)$           | $35.8 (\pm 0.9)$            | $44.2 (\pm 0.8)$    | $40.1 (\pm 2.5)$  |
| DRNet       | $72.8(\pm 1.3)$   | $71.3(\pm 2.5)$   | $71.9(\pm 1.8)$            | $65.5 (\pm 0.7)$    | $66.8 (\pm 1.8)$           | $70.4 (\pm 2.7)$            | $85.4 (\pm 0.7)$    | $77.6 (\pm 1.4)$  |
| MRNet       | $34.4 (\pm 3.4)$  | $45.6(\pm 6.1)$   | $33.1 (\pm 8.9)$           | $30.1 (\pm 9.1)$    | $24.1 (\pm 1.8)$           | $23.7 (\pm 4.5)$            | $49.6 (\pm 1.5)$    | $35.2 (\pm 1.1)$  |
| PrAE        | $62.3 (\pm 0.9)$  | $73.3(\pm 2.0)$   | $\underline{75.0}$ (± 3.3) | $41.4 (\pm 5.9)$    | $83.0 (\pm 0.7)$           | 82.9 (± 1.2)                | $47.1 (\pm 4.7)$    | $33.0 (\pm 1.7)$  |
| PredRNet    | $62.3 (\pm 2.2)$  | $64.2 (\pm 1.2)$  | $61.6 (\pm 5.8)$           | $52.0 (\pm 5.2)$    | $64.3 (\pm 1.5)$           | $64.1 (\pm 0.8)$            | $75.3 (\pm 0.9)$    | $54.8 (\pm 6.1)$  |
| RelBase     | $73.0(\pm 1.8)$   | $73.8(\pm 4.0)$   | $72.4 (\pm 0.6)$           | $65.1 (\pm 0.2)$    | $69.7 (\pm 2.9)$           | $\underline{73.0}(\pm 5.1)$ | $83.1 (\pm 2.7)$    | $73.8(\pm 2.5)$   |
| SCL         | $75.6 (\pm 10.1)$ | $84.0 (\pm 6.3)$  | $74.8 (\pm 10.3)$          | $66.4 (\pm 14.7)$   | $73.2 (\pm 13.1)$          | $71.6 (\pm 12.4)$           | $81.1 (\pm 7.5)$    | $77.6 (\pm 8.8)$  |
| SRAN        | $42.1 (\pm 2.3)$  | $47.7 (\pm 4.3)$  | $35.5 (\pm 0.9)$           | $31.4 (\pm 0.7)$    | $41.8 (\pm 4.4)$           | $42.2 (\pm 4.1)$            | $51.3(\pm 1.9)$     | $44.6 (\pm 1.2)$  |
| STSN        | $39.9(\pm 14.7)$  | $47.1 (\pm 9.6)$  | $35.8 (\pm 11.4)$          | $32.3 (\pm 10.5)$   | $38.5 (\pm 14.2)$          | $39.5 (\pm 14.2)$           | $49.1 \ (\pm 26.5)$ | $38.1 (\pm 18.6)$ |
| WReN        | $18.0 (\pm 0.4)$  | $20.7 (\pm 1.0)$  | $18.8 (\pm 0.3)$           | $17.4 (\pm 0.7)$    | $15.8 (\pm 0.6)$           | $15.1 (\pm 0.5)$            | $19.3 (\pm 0.7)$    | $18.8 (\pm 0.4)$  |
| PoNG (ours) | $81.4 (\pm 3.1)$  | $88.3 (\pm 6.2)$  | <b>86.5</b> $(\pm 5.0)$    | <b>79.8</b> (± 3.3) | $\underline{73.4}$ (± 3.0) | $71.1~(\pm 2.1)$            | 87.2 $(\pm 2.8)$    | 83.5 (±1.0)       |
|             | •                 |                   |                            |                     |                            |                             |                     |                   |

Table 33: A/Color-Arithmetic. The table presents results from Table 10 extended to each matrix configuration. 

|             | Mean              | Center            | 2x2-Grid                | 3x3-Grid          | L-R               | U-D               | O-IC                       | O-IG                       |
|-------------|-------------------|-------------------|-------------------------|-------------------|-------------------|-------------------|----------------------------|----------------------------|
| ALANS       | $18.3 (\pm 6.6)$  | $18.6 (\pm 6.4)$  | $18.2 (\pm 5.5)$        | $17.0 (\pm 5.0)$  | $20.0 (\pm 8.2)$  | $18.8 (\pm 7.6)$  | $18.1 (\pm 8.0)$           | $17.4 (\pm 5.7)$           |
| CPCNet      | $45.9(\pm 2.7)$   | $41.2(\pm 3.3)$   | $39.2 (\pm 5.7)$        | $35.6 (\pm 5.2)$  | $50.5 (\pm 1.0)$  | $47.7 (\pm 0.8)$  | $62.2 (\pm 1.7)$           | $45.5 (\pm 2.7)$           |
| CNN-LSTM    | $17.1 (\pm 3.7)$  | $18.8 (\pm 5.9)$  | $17.3 (\pm 4.2)$        | $17.3 (\pm 3.4)$  | $15.8 (\pm 3.2)$  | $15.6 (\pm 2.4)$  | $17.2 (\pm 3.9)$           | $17.8 (\pm 3.6)$           |
| CoPINet     | $35.2 (\pm 0.5)$  | $35.9(\pm 0.5)$   | $31.4 (\pm 0.5)$        | $27.7 (\pm 0.7)$  | $34.1 (\pm 1.2)$  | $33.5 (\pm 1.6)$  | $43.4 (\pm 0.7)$           | $40.3 (\pm 0.5)$           |
| DRNet       | $66.7 (\pm 1.2)$  | $60.0(\pm 2.2)$   | $69.1 (\pm 1.8)$        | $62.5 (\pm 1.9)$  | $63.3 (\pm 2.5)$  | 63.2 (± 4.2)      | 77.5 (± 3.0)               | $71.6 (\pm 2.2)$           |
| MRNet       | $35.7 (\pm 5.9)$  | $38.9 (\pm 14.1)$ | $38.8(\pm 3.1)$         | $32.4 (\pm 3.2)$  | $20.6 (\pm 7.5)$  | $20.4 (\pm 6.9)$  | $46.6 (\pm 7.0)$           | $52.4 (\pm 0.8)$           |
| PrAE        | $43.0 (\pm 26.5)$ | $47.2 (\pm 30.6)$ | $52.5~(\pm 34.9)$       | $32.0\ (\pm16.5)$ | $54.6 (\pm 36.0)$ | $53.4~(\pm35.3)$  | $34.4 \ (\pm 18.8)$        | $26.8 \ (\pm 13.8)$        |
| PredRNet    | $56.9(\pm 1.4)$   | $50.8 (\pm 2.2)$  | $62.5 (\pm 1.8)$        | $54.5 (\pm 1.9)$  | $55.3 (\pm 0.4)$  | $53.7 (\pm 1.5)$  | $68.8 (\pm 0.2)$           | $53.4 (\pm 6.0)$           |
| RelBase     | $66.2 (\pm 1.0)$  | $60.4(\pm 1.6)$   | $69.4 (\pm 0.4)$        | $61.9 (\pm 3.0)$  | $63.1 (\pm 1.1)$  | $62.4 (\pm 0.4)$  | $74.4(\pm 1.9)$            | $\underline{72.0}$ (± 1.8) |
| SCL         | $60.0 (\pm 4.1)$  | $62.8 (\pm 5.4)$  | $58.5 (\pm 2.9)$        | $54.0 (\pm 2.2)$  | $54.6 (\pm 6.6)$  | $52.3 (\pm 5.8)$  | $69.1 (\pm 5.2)$           | $68.8 (\pm 1.4)$           |
| SRAN        | $39.9(\pm 2.7)$   | $39.8 (\pm 6.0)$  | $37.6 (\pm 2.7)$        | $32.0 (\pm 2.1)$  | $39.8 (\pm 3.1)$  | $38.6 (\pm 4.0)$  | $49.3 (\pm 1.3)$           | $42.6 (\pm 1.3)$           |
| STSN        | $25.7 (\pm 10.6)$ | $28.5 (\pm 7.2)$  | $26.1 (\pm 8.9)$        | $23.3 (\pm 8.7)$  | $22.7 (\pm 11.0)$ | $23.5\;(\pm10.3)$ | $30.4~(\pm 17.4)$          | $25.9 \ (\pm 11.5)$        |
| WReN        | $17.1 (\pm 0.2)$  | $19.1 (\pm 0.8)$  | $16.6 (\pm 0.6)$        | $15.3 (\pm 0.2)$  | $16.1 (\pm 0.8)$  | $16.0 (\pm 0.4)$  | $18.7 (\pm 0.7)$           | $18.3 (\pm 0.6)$           |
| PoNG (ours) | $70.0 (\pm 4.1)$  | 64.4 (± 7.9)      | <b>74.4</b> $(\pm 3.9)$ | $69.8 (\pm 3.2)$  | $63.4 (\pm 3.3)$  | $62.1~(\pm 5.1)$  | $\underline{76.0}$ (± 4.5) | <b>80.1</b> (± 2.1)        |

Table 34: A/Color-DistributeThree. The table presents results from Table 10 extended to each matrix configuration. 

| 1390    |             |                    |                             |                     |                     |                     |                     |                            |                              |
|---------|-------------|--------------------|-----------------------------|---------------------|---------------------|---------------------|---------------------|----------------------------|------------------------------|
| 1391    |             | Mean               | Center                      | 2x2-Grid            | 3x3-Grid            | L-R                 | U-D                 | O-IC                       | O-IG                         |
| 1392    | ALANS       | $22.4 (\pm 7.7)$   | $20.6 (\pm 7.3)$            | $22.0 (\pm 5.2)$    | $23.2 (\pm 8.1)$    | $23.8 (\pm 9.7)$    | $24.2 (\pm 9.8)$    | $23.9 (\pm 9.8)$           | $19.0 (\pm 4.3)$             |
| 1393    | CPCNet      | $37.8 \ (\pm 0.9)$ | $31.0\;(\pm1.8)$            | $28.7 \ (\pm 2.9)$  | $28.6~(\pm 2.4)$    | $41.7\;(\pm0.0)$    | $41.0 \; (\pm 1.3)$ | $54.0 \ (\pm 0.8)$         | $39.9 \ (\pm 0.3)$           |
| 1000    | CNN-LSTM    | $20.6 (\pm 6.7)$   | $24.6 \ (\pm 10.1)$         | $20.1 (\pm 4.9)$    | $19.2 (\pm 5.2)$    | $18.7 (\pm 5.3)$    | $17.9 (\pm 4.1)$    | $22.5 (\pm 9.2)$           | $21.7 (\pm 7.9)$             |
| 1394    | CoPINet     | $26.9 (\pm 0.5)$   | $25.8 (\pm 1.0)$            | $24.1 (\pm 0.1)$    | $21.8 (\pm 0.7)$    | $22.9 (\pm 1.6)$    | $23.4 (\pm 0.2)$    | $34.4 (\pm 1.2)$           | $35.6 (\pm 1.0)$             |
| 1395    | DRNet       | $63.2 (\pm 0.3)$   | $57.1 (\pm 2.0)$            | $64.1 (\pm 1.7)$    | $58.2(\pm 1.8)$     | $61.1 (\pm 0.9)$    | $61.2 (\pm 2.4)$    | $72.9(\pm 2.0)$            | $67.9 (\pm 1.6)$             |
| 1396    | MRNet       | $18.6 (\pm 0.1)$   | $20.9(\pm 3.4)$             | $27.7 (\pm 3.3)$    | $24.0 \ (\pm 6.6)$  | $11.9 (\pm 6.1)$    | $13.2 (\pm 2.5)$    | $14.3 (\pm 4.6)$           | $18.3 (\pm 6.6)$             |
|         | PrAE        | $55.1 (\pm 0.8)$   | $61.9(\pm 1.4)$             | $63.9 (\pm 2.7)$    | $41.8 (\pm 1.5)$    | $69.6 (\pm 0.9)$    | $69.5 (\pm 0.9)$    | $45.7 (\pm 2.1)$           | $33.5(\pm 3.9)$              |
| 1397    | PredRNet    | $48.5 (\pm 0.9)$   | $38.7 (\pm 0.5)$            | $51.2 (\pm 3.8)$    | $43.7 (\pm 4.5)$    | $46.2 (\pm 0.3)$    | $45.8 (\pm 0.3)$    | $61.0 (\pm 0.8)$           | $52.7 (\pm 3.4)$             |
| 1398    | RelBase     | $65.7 (\pm 4.6)$   | $60.8 (\pm 6.2)$            | $64.9(\pm 3.8)$     | $57.9 (\pm 5.5)$    | $67.6 (\pm 5.3)$    | $66.1 (\pm 4.8)$    | $\underline{75.6}$ (± 2.8) | $67.2 (\pm 4.8)$             |
| 1399    | SCL         | $63.9 \ (\pm 4.3)$ | $\underline{70.8}(\pm 4.5)$ | $60.3 \ (\pm 4.3)$  | $54.5~(\pm 3.7)$    | $59.7~(\pm 5.8)$    | $57.8~(\pm 6.4)$    | $72.0 \ (\pm 4.0)$         | $\underline{71.8} (\pm 2.2)$ |
| 4 4 9 9 | SRAN        | $34.6 (\pm 3.6)$   | $38.4 (\pm 8.4)$            | $32.5 (\pm 4.5)$    | $30.1 (\pm 2.8)$    | $32.6 (\pm 3.1)$    | $31.5 (\pm 3.1)$    | $40.8 (\pm 2.3)$           | $36.6 (\pm 2.2)$             |
| 1400    | STSN        | $20.7 (\pm 7.7)$   | $20.8 (\pm 6.4)$            | $19.0 (\pm 6.8)$    | $19.1 (\pm 7.7)$    | $18.2 (\pm 6.5)$    | $16.1 (\pm 4.7)$    | $27.5 (\pm 14.1)$          | $25.0 (\pm 10.6)$            |
| 1401    | WReN        | $17.7 (\pm 0.6)$   | $21.0 (\pm 0.8)$            | $17.3 (\pm 0.5)$    | $16.4 (\pm 0.4)$    | $16.0 (\pm 0.7)$    | $16.3 (\pm 1.1)$    | $18.0 (\pm 1.7)$           | $19.1 (\pm 2.1)$             |
| 1402    | PoNG (ours) | 81.3 $(\pm 1.6)$   | <b>84.1</b> (± 4.9)         | <b>82.7</b> (± 2.1) | <b>77.6</b> (± 1.4) | <b>80.6</b> (± 2.8) | <b>80.9</b> (± 3.5) | $84.4 (\pm 2.1)$           | <b>79.1</b> $(\pm 3.5)$      |



# 1404 F DATASHEETS FOR DATASETS

In what follows, we provide the description of the introduced datasets following the Datasheets for Datasets template (Gebru et al., 2021).

For what purpose was the dataset created?
Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.

Motivation

The datasets were created to study generalization and knowledge transfer abilities of AVR models.

# Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?

1423 1424 Hidden for blind review.

Who funded the creation of the dataset? If
there is an associated grant, please provide the
name of the grantor and the grant name and
number.

1429 1430 Hidden for blind review.

#### Any other comments?

1433 None.

1434 1435

1436

1455

1432

1409 1410

1411

1419

#### Composition

What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and interactions between them; nodes and edges)? Please provide a description.

Each dataset instance represents a single Raven's
Progressive Matrix, which is a typical task used in human IQ tests.

How many instances are there in total (of each type, if appropriate)?

Each regime in Attributeless-I-RAVEN as well as the I-RAVEN-Mesh dataset contains 70 000 instances. The training, validation, and test splits contain 42 000, 14 000, and 14 000 matrices, resp. All together there are 770 000 (11 × 70 000) instances.

1456Does the dataset contain all possible in-<br/>stances or is it a sample (not necessarily ran-<br/>dom) of instances from a larger set?

the dataset is a sample, then what is the larger set? Is the sample representative of the larger set (e.g., geographic coverage)? If so, please describe how this representativeness was validated/verified. If it is not representative of the larger set, please describe why not (e.g., to cover a more diverse range of instances, because instances were withheld or unavailable).

The datasets contain a fixed number of instances generated with the data generator. Using a fixed seed ensures reproducibility of the generation process. The data generator allows to configure the number of generated samples.

What data does each instance consist of? "Raw" data (e.g., unprocessed text or images) or features? In either case, please provide a description.

Each RPM instance comprises 16 images that represent the RPM panels, a corresponding index of the correct answer and a representation of rules that govern the matrix. Section 3 of the paper provides additional details. Each instance is packaged as a separate file in the NPZ format, which is a widely-used binary format to store compressed NumPy arrays.

Is there a label or target associated with each instance? If so, please provide a description.

#### See above.

**Is any information missing from individual instances?** If so, please provide a description, explaining why this information is missing (e.g., because it was unavailable). This does not include intentionally removed information, but might include, e.g., redacted text.

There's no missing data.

Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)? If so, please describe how these relationships are made explicit.

There are no relationships between individual instances.

Are there recommended data splits (e.g., training, development/validation, testing)? If so, please provide a description of these splits, explaining the rationale behind them.

1458 The datasets are split into training, validation 1459 and test splits. Each generated instance con-1460 tains the split name in its filename, e.g., the 1461 RAVEN\_1\_train.npz file belongs to the train 1462 split.

1464 Are there any errors, sources of noise, or re-1465 dundancies in the dataset? If so, please provide a description. 1466

1467 To the best of our knowledge there are no er-1468 rors, sources of noise, nor redundancies in the 1469 datasets.

1470

1463

1471 Is the dataset self-contained, or does it link to 1472 or otherwise rely on external resources (e.g., 1473 websites, tweets, other datasets)? If it links to 1474 or relies on external resources, a) are there quar-1475 antees that they will exist, and remain constant, over time; b) are there official archival versions 1476 of the complete dataset (i.e., including the ex-1477 ternal resources as they existed at the time the 1478 dataset was created); c) are there any restric-1479 tions (e.g., licenses, fees) associated with any of 1480 the external resources that might apply to a fu-1481 ture user? Please provide descriptions of all ex-1482 ternal resources and any restrictions associated 1483 with them, as well as links or other access points, 1484 as appropriate.

1485 Both datasets are self-contained. 1486

1487

Does the dataset contain data that might be 1488 considered confidential (e.g., data that is pro-1489 tected by legal privilege or by doctor-patient 1490 confidentiality, data that includes the content 1491 of individuals non-public communications)? 1492 If so, please provide a description. 1493

No. 1494

1495 Does the dataset contain data that, if viewed 1496 directly, might be offensive, insulting, threat-1497 ening, or might otherwise cause anxiety? If 1498 so, please describe why.

1500 No.

1501

1499

1502 Does the dataset relate to people? If not, you 1503 may skip the remaining questions in this section. 1504

No. 1505

1506 Does the dataset identify any subpopulations 1507 (e.g., by age, gender)? If so, please describe 1508 how these subpopulations are identified and pro-1509 vide a description of their respective distributions 1510 within the dataset. 1511

N/A.

Is it possible to identify individuals (i.e., one or more natural persons), either directly or indirectly (i.e., in combination with other data) from the dataset? If so, please describe how.

N/A.

Does the dataset contain data that might be considered sensitive in any way (e.g., data that reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or union memberships, or locations; financial or health data; biometric or genetic data; forms of government identification, such as social security numbers; criminal history)? If so, please provide a description.

N/A.

Any other comments?

None.

#### **Collection Process**

How was the data associated with each instance acquired? Was the data directly observable (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-ofspeech tags, model-based guesses for age or language)? If data was reported by subjects or indirectly inferred/derived from other data, was the data validated/verified? If so, please describe how.

The data was generated with a computer program.

What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sensor, manual human curation, software program, software API)? How were these mechanisms or procedures validated?

We extended the data generation code used to create I-RAVEN: https://github.com/ husheng12345/SRAN. A subset of the dataset was reviewed manually to ensure correctness of the generated matrices.

#### If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic, probabilistic with specific sampling probabilities)?

The dataset is produced by a generator that creates new RPM instances subject to specified constraints through a pseudo-random process. We use a fixed seed to ensure reproducibility of the generation process.

1512 Who was involved in the data collection pro-1513 cess (e.g., students, crowdworkers, contrac-1514 tors) and how were they compensated (e.g., (if appropriate). 1515 how much were crowdworkers paid)? N/A. 1516 The data generator has been written by the au-1517 thors of this paper without delegating the work 1518 to other individuals. 1519 1520 Over what timeframe was the data collected? 1521 Does this timeframe match the creation time-1522 frame of the data associated with the inumentation. 1523 stances (e.g., recent crawl of old news arti-1524 cles)? If not, please describe the timeframe in N/A. which the data associated with the instances was 1525 created. 1526 1527 Development of the datasets lasted from January 1528 2022 to September 2024. rights. 1529 Were any ethical review processes conducted 1530 (e.g., by an institutional review board)? If so, 1531 please provide a description of these review pro-1532 cesses, including the outcomes, as well as a link 1533 or other access point to any supporting documen-1534 tation. 1535 N/A. 1536 1537 Does the dataset relate to people? If not, you 1538 may skip the remaining questions in this section. 1539 1540 No. 1541 1542 Did you collect the data from the individuals 1543 in question directly, or obtain it via third parties or other sources (e.g., websites)? 1544 1545 N/A. 1546 1547 Were the individuals in question notified "raw" data. 1548 about the data collection? If so, please de-1549 scribe (or show with screenshots or other infor-N/A. mation) how notice was provided, and provide a 1550 link or other access point to, or otherwise repro-1551 duce, the exact language of the notification itself. 1552 1553 N/A. 1554 1555 Did the individuals in question consent to the 1556 collection and use of their data? If so, please describe (or show with screenshots or other in-1557 formation) how consent was requested and pro-1558 vided, and provide a link or other access point 1559 ing the datasets. to, or otherwise reproduce, the exact language to 1560

#### Any other comments?

None.

If consent was obtained, were the consenting 1564 individuals provided with a mechanism to re-1565 voke their consent in the future or for certain

which the individuals consented.

1561

1562

1563

N/A.

Uses

Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted? If so, please provide a description of this analysis, including the outcomes, as well as a link or other access point to any supporting doc-

#### Any other comments?

We bear all responsibility in case of violation of

#### Preprocessing/cleaning/labeling

Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)? If so, please provide a description. If not, you may skip the remainder of the questions in this section.

The data created by the generator is ready to be used in a model. No preprocessing, cleaning, or labeling is required.

Was the "raw" data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)? If so, please provide a link or other access point to the

Is the software used to preprocess/clean/label the instances available? If so, please provide a link or other access point.

No specific software is required to preprocess, clean, or label the instances. The released code repository contains the code required to reproduce all experiments from the paper, which can be used as a reference implementation for loadHas the dataset been used for any tasks already? If so, please provide a description.

1569 The datasets have been used to conduct experiments presented in the paper.

1571
1572
1573
1573
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574
1574</l

1575 N/A.

#### 1576 What (other) tasks could the dataset be used for?

1579 The datasets could be used in a multi-task setting to improve abstract reasoning capabilities of computer vision models.

1582 Is there anything about the composition of the dataset or the way it was collected and pre-1584 processed/cleaned/labeled that might impact 1585 future uses? For example, is there anything that a future user might need to know to avoid uses 1587 that could result in unfair treatment of individuals or groups (e.g., stereotyping, quality of service is-1589 sues) or other undesirable harms (e.g., financial harms, legal risks) If so, please provide a descrip-1590 tion. Is there anything a future user could do to 1591 mitigate these undesirable harms? 1592

We do not see any undesirable harms that could apply to future users of the datasets.

1596Are there tasks for which the dataset should1597not be used? If so, please provide a description.

The datasets should not be used in human IQ
tests, as they were explicitly designed to assess
generalization and knowledge transfer abilities
of deep learning models.

Any other comments?

5 None.

1606

1604

1595

160

1608

1609 1610 How will the dataset will be distributed (e.g., tarball on website, API, GitHub) Does the dataset have a digital object identifier (DOI)?

The datasets will be shared on GitHub.

#### When will the dataset be distributed?

The datasets will become publicly available after paper acceptance.

Will the dataset be distributed under a copyright or other intellectual property (IP) license, and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or ToU, as well as any fees associated with these restrictions.

The code repository is released under the GPL-3.0 license. This follows the license associated with the generators of the base datasets – RAVEN (https://github.com/WellyZhang/RAVEN) and I-RAVEN (https://github.com/husheng12345/SRAN). The datasets introduced in this paper are released under the CC license.

Have any third parties imposed IP-based or other restrictions on the data associated with the instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these restrictions.

No.

Do any export controls or other regulatory restrictions apply to the dataset or to individual instances? If so, please describe these restrictions, and provide a link or other access point to, or otherwise reproduce, any supporting documentation.

No.

#### Any other comments?

None.

Maintenance

# Who will be supporting/hosting/maintaining the dataset?

Hidden for blind review.

1611 ties outside of the entity (e.g., company, in1612 stitution, organization) on behalf of which the N
1613 dataset was created? If so, please provide a
1614 description.
1615 The datasets will become publicly available after

Distribution

Will the dataset be distributed to third par-

The datasets will become publicly available after
paper acceptance. Additionally, as discussed in
Section 4 ("Reproducibility"), the attached code
allows for generation of all datasets from scratch,
eliminating the dependency on file-hosting services required to distribute the data.

#### 1620 How can the owner/curator/manager of the N/A. 1621 dataset be contacted (e.g., email address)? 1622

Hidden for blind review. 1623

1624 Is there an erratum? If so, please provide a link 1625 or other access point. 1626

No. 1627

1628

1636

Will the dataset be updated (e.g., to correct 1629 labeling errors, add new instances, delete in-1630 stances)? If so, please describe how often, by 1631 whom, and how updates will be communicated 1632 to users (e.g., mailing list, GitHub)? 1633

Future changes will be documented in release 1634 notes in the code repository. 1635

If the dataset relates to people, are there ap-1637 plicable limits on the retention of the data as-1638 sociated with the instances (e.g., were indi-1620 duals in question told that their data would retained for a fixed period of time and then leted)? If so, please describe these limits and plain how they will be enforced.

Will older versions of the dataset continue to be supported/hosted/maintained? If so, please describe how. If not, please describe how its obsolescence will be communicated to users.

Older versions of the dataset will be available in the history of the code repository.

If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for them to do so? If so, please provide a description. Will these contributions be validated/verified? If so, please describe how. If not, why not? Is there a process for communicating/distributing these contributions to other users? If so, please provide a description.

Contributions are welcome. GitHub Issues of the code repository will be used to communicate between contributors and project maintainers.

#### Any other comments?

None.

| 1039 | vio |
|------|-----|
| 1640 | be  |
| 1641 | de  |
| 1642 | ex  |
| 1643 |     |
| 1644 |     |
| 1645 |     |
| 1646 |     |
| 1647 |     |
| 1648 |     |
| 1649 |     |
| 1650 |     |
| 1651 |     |
| 1652 |     |
| 1653 |     |
| 1654 |     |
| 1655 |     |
| 1656 |     |
| 1657 |     |
| 1658 |     |
| 1659 |     |
| 1660 |     |
| 1661 |     |
| 1662 |     |
| 1663 |     |
| 1664 |     |
| 1665 |     |
| 1666 |     |
| 1667 |     |
| 1668 |     |
| 1669 |     |
| 1670 |     |
| 1671 |     |
| 1672 |     |
| 1673 |     |