

326 A Models

327 Figure 1 shows an illustration of how we use transformers for sequence classification (left) and
 328 sequence transduction (right) problems. In the classification setting the input is a sequence $(x_1 \cdots x_n)$
 329 and the output is a discrete label y . In the transduction setting, the input $(x_1 \cdots x_n)$ and output
 330 $(y_1 \cdots y_m)$ are sequences. z is an embedding vector we refer to as the *task embedding* and appears in
 331 the input to the transformer, in addition to the input sequence. The task embedding z is task specific,
 332 and is inferred on the fly for each task during training. Learning a new task \mathcal{T} at test time involves
 333 inferring the corresponding task embedding $z_{\mathcal{T}}$, leaving the rest of the model parameters untouched.

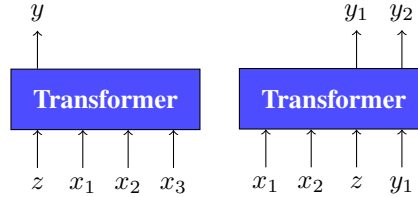


Figure 1: Illustration of how we use transformers for sequence classification (left) and sequence transduction (right) problems.

334 B Hyperparameters

335 Unless otherwise specified our transformer has 4 layers with an embedding size of 128. We use the
 336 Adam optimizer (Kingma & Ba, 2015) for both outer and inner loop optimization. The maximum
 337 number of task embedding optimization steps is set to 25 during training. We train a single model
 338 using N samples at training time, treating N as a hyperparameter and apply it to k -shot problems
 339 with different values of k at test time. The size of the task embedding was set to match the embedding
 340 dimension of the transformer (128).

341 C Baselines

342 **Matching Networks** (Vinyals et al., 2016) This model computes the similarity between the query
 343 input and each training example. In our implementation we adapted the original vision model
 344 to our sequence learning setting by feeding the concatenation of both sequences to a transformer
 345 which outputs a single scalar score, which can be interpreted as the similarity between the two input
 346 sequences. The prediction is a convex sum of the training example labels, the weights being the
 347 similarity scores. We consider Matching Networks only in our classification setting, as it is not
 348 straightforward to use it for transduction. Matching networks are learned using the multi-task training
 349 loss; at test time they are applied to new tasks without any finetuning.

350 **SNAIL** (Mishra et al., 2018) This model is similar to the task-agnostic transformer except that the
 351 input is augmented with the concatenation of all input-output training pairs. Similarly to Matching
 352 Networks, SNAIL is learned using the multi-task training loss and applied to new test tasks without
 353 finetuning. For both Matching Networks and SNAIL, we construct training episodes by sampling k
 354 training examples to define a task, to match with the test scenario. We thus train different models for
 355 each k -shot problem.

356 **MAML** (Finn et al., 2017) All model parameters are trained using MAML, with the same model
 357 architecture as TAM. The entire model is fine-tuned on test tasks.

358 **CAVIA** (Zintgraf et al., 2019) Similar to TAM, CAVIA has a set of task-specific parameters and
 359 shared parameters. The training algorithm is similar to MAML, but inner loop updates are performed
 360 on the task-specific parameters as opposed to the entire model. Our implementation of CAVIA,
 361 MAML uses the *higher* library (Grefenstette et al., 2019).

D Sequence transformations used to construct classification and transduction tasks

In tables 5, 6 we describe the transformations used to construct classification and transduction tasks, respectively.

	Transformation	Description
S_1	mul v	Elementwise multiply by v
	add v	Elementwise add v
	div v	Elementwise integer division by v
	mod v	Elementwise modulo v operation
S_2	(not) multiple of v	Extract subset of integers that are (not) multiples of v
	(not) greater of v	Extract subset of integers that are (not) greater than v
	(do not) have exactly v divisors	Extract subset of integers that (do not) have exactly v divisors
S_3	count	Sequence length
	min	Smallest integer in sequence
	max	Largest integer in sequence
	mean	Mean of sequence elements
	median	Median of sequence elements
	mode	Mode of sequence elements
	first	First element in sequence
	last	Last element in sequence
	max-min	Difference between largest and smallest elements in sequence
	middle	Element in the middle position of sequence

Table 5: Sequence transformations used to construct classification tasks and their descriptions. Each transformation takes a sequence as input and outputs a sequence (transformations in S_1 and S_2), or a single integer (transformations in S_3).

	Transformation	Description
S_1	mul v	Elementwise multiply by v
	add v	Elementwise add v
	div v	Elementwise integer division by v
	mod v	Elementwise modulo v operation
S_2	reverse v with v'	Replace all occurrences of v in the sequence with v'
	replace x_i with $f(x_i, x_j)$	Replace element x_i with one of the following: $\{ax_i + b, x_j, \text{abs}(x_i - x_j), x_i + x_j\}$ where a, b are integer constants and x_i, x_j are elements of the sequence at position i, j respectively
S_3	sort ascending	Sort the sequence in ascending order
	sort descending	Sort the sequence in descending order
	reverse	Reverse the sequence
	swap(x_i, x_j)	Swap elements at positions i, j of the sequence
	shift right v	Cyclic shift the sequence right by v positions

Table 6: Sequence transformations used to construct transduction tasks and their descriptions. Each transformation takes a sequence as input and outputs a sequence.

366 E Compositional TAM

367 Algorithm 2 presents the training algorithm for compositional TAM. We draw a training task $\mathcal{T}^{\text{train}}$
 368 with primitive ids $T_1 = i_1, T_2 = i_2, T_3 = i_3$ respectively in line 3. These primitive ids index into the
 369 primitive embedding table θ_e . We pretend that one of the primitives is unknown, and to illustrate
 370 the algorithm, we assume without loss of generality that $T_2 = i_2$ is unknown (line 5). In the inner
 371 loop optimization, we infer an embedding z for this unknown primitive using gradient descent, while
 372 using the primitive embedding table to load the known primitive embeddings ($\theta_e[i_1], \theta_e[i_3]$ in this
 373 case (lines 8, 9)).

Algorithm 2: Compositional TAM for k-shot Learning

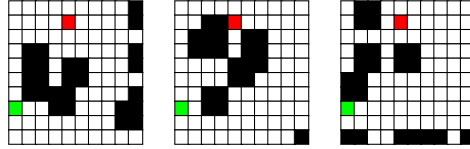
Input : Training tasks $\mathcal{T}_1^{\text{train}}, \dots, \mathcal{T}_N^{\text{train}}$
Output : Model parameters θ , primitive embeddings θ_e

- 1 $\theta' = \theta \cup \theta_e$;
- 2 **repeat**
- 3 Sample training task $\mathcal{T}^{\text{train}}$ with primitive ids $T_1 = i_1, T_2 = i_2, T_3 = i_3$;
- 4 Sample k training examples from the task $\{(x^j, y^j)_{j=1, \dots, k}\} \sim \mathcal{T}^{\text{train}}$;
- 5 Pretend one of the primitives (chosen at random) is unknown, say T_2
- 6 Initialize $z = 0, \Delta\theta' = 0$;
- 7 **while** *loss improves and max iterations not reached* **do**
- 8 $z \leftarrow z - \nabla_z \sum_{j=1}^k -\log p(y^j | x^j, z_1 = \theta_e[i_1], z_2 = z, z_3 = \theta_e[i_3]; \theta')$
- 9 $\Delta\theta' \leftarrow \Delta\theta' - \nabla_{\theta'} \sum_{j=1}^k -\log p(y^j | x^j, z_1 = \theta_e[i_1], z_2 = z, z_3 = \theta_e[i_3]; \theta')$
- 10 $\theta' \leftarrow \theta' + \Delta\theta'$
- 11 **until** *max training iterations*;

F Path-finding task

F.1 Non-compositional path-finding task

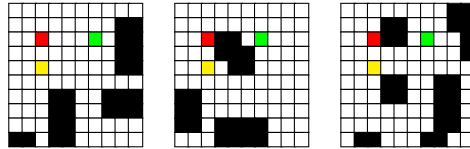
We present an example task from the path-finding task below. The grids are 10×10 . The following task is defined by the start position (7, 0) and end position (1, 4), indicated by the green and red squares, respectively. Each example in the task corresponds to a particular configuration of obstacles in the grid. The source sequence represents the locations of obstacles. The obstacles are represented by the top left position of a 2×2 blob. The target sequence represents the optimal path from source to target. Source and target sequences consist of rasterized grid coordinates (Eg. rasterized start and end positions are 70 and 14, respectively). In addition, elements of the target sequence have an offset of 100 (Eg. rasterized position 14 is represented as 114).



- Source: [39, 78, 51, 9, 31, 63, 44, 69], Target: [170, 160, 150, 140, 130, 121, 112, 103, 114]
- Source: [12, 35, 99, 22, 62, 44, 25, 21], Target: [170, 161, 152, 143, 134, 124, 114]
- Source: [90, 99, 1, 96, 34, 50, 94, 31], Target: [170, 171, 162, 152, 143, 133, 123, 114]

F.2 Compositional path-finding task

In the compositional setting, we require the optimal path to pass through a way-point, indicated in yellow in the following grids. A task is thus defined by a start position, end position and way-point position. The possible values for each of these three parameters represent the primitives in this compositional setting.



- Source: [63, 38, 90, 93, 73, 68, 18, 67], Target: [126, 115, 124, 133, 142, 131, 122]
- Source: [95, 60, 95, 70, 23, 34, 83, 85], Target: [126, 115, 104, 113, 122, 131, 142, 131, 122]
- Source: [91, 29, 57, 96, 8, 53, 77, 13], Target: [126, 125, 134, 133, 142, 131, 122]

397 G Results with error bars

398 Tables 7, 8 show the error bars on results in the main paper.

Model	Sequence Classification				Sequence Transduction				Path Finding			
	1	5	10	20	1	5	10	20	1	5	10	20
Task Agnostic	40.50	64.75	74.25	82.50	6.27	5.26	4.71	4.01	3.17	1.75	1.55	1.39
	± 1.73	± 1.29	± 1.29	± 1.58	± 0.17	± 0.04	± 0.05	± 0.07	± 0.27	± 0.03	± 0.02	± 0.01
Multitask	38.75	66.00	77.50	87.50	13.80	6.80	5.18	2.91	6.39	1.98	1.64	1.44
	± 0.96	± 0.82	± 1.29	± 0.58	± 3.36	± 0.26	± 1.01	± 0.49	± 1.96	± 0.12	± 0.05	± 0.02
Matching	43.00	58.75	64.50	67.00	—				—			
	± 1.15	± 2.45	± 1.83	± 1.41								
SNAIL	43.00	44.00	68.25	67.50	2.48	3.80	4.98	4.11	2.63	1.95	3.47	3.04
	± 1.41	± 2.00	± 1.26	± 4.43	± 0.38	± 0.38	± 0.03	± 2.94	± 0.27	± 0.25	± 1.56	± 1.01
MAML	39.60	63.40	71.80	78.80	6.73	5.84	5.19	4.20	5.49	2.05	1.65	1.44
	± 0.55	± 0.89	± 0.84	± 0.84	± 0.16	± 0.2	± 0.08	± 0.10	± 0.85	± 0.03	± 0.01	± 0.01
CAVIA	43.00	78.00	87.00	91.00	11.05	2.75	1.78	1.53	2.21	1.31	1.25	1.21
	± 0.58	± 1.26	± 0.50	± 0.58	± 2.94	± 0.51	± 0.14	± 0.06	± 0.21	± 0.03	± 0.03	± 0.02
TAM	40.50	75.50	89.50	94.50	8.47	2.92	1.47	1.15	1.82	1.27	1.22	1.17
	± 0.82	± 0.50	± 0.58	± 0.82	± 1.42	± 0.67	± 0.18	± 0.03	± 0.08	± 0.01	± 0.01	± 0.01

Table 7: k -shot sequence classification and sequence transduction experiments on our three benchmarks for $k \in \{1, 5, 10, 20\}$. The metric for sequence classification is average accuracy on test tasks (higher is better). On the transduction tasks, the performance metric is average perplexity on test tasks (lower is better). Random performance is at 25% accuracy (classification) and 12 perplexity points (other two tasks). Entries in smaller font are error bars, and they are estimated on 4 trials varying the model initialization.

Model	Sequence Classification				Sequence Transduction				Path Finding			
	1	5	10	20	1	5	10	20	1	5	10	20
Multitask	57.50	74.00	81.00	88.5	43.32	7.50	3.48	2.16	3.08	1.62	1.32	1.23
	± 3.51	± 5.48	± 4.24	± 2.08	± 10.87	± 0.43	± 0.12	± 0.05	± 0.61	± 0.33	± 0.05	± 0.02
Matching	61.25	69.50	72.25	67.5	—				—			
	± 4.35	± 5.45	± 6.08	± 5.92								
SNAIL	63.5	71.75	76.25	80.25	7.00	6.24	6.93	17.10	1.53	1.71	3.36	4.22
	± 4.80	± 4.86	± 2.75	± 2.87	± 2.03	± 0.27	± 3.25	± 9.66	± 0.29	± 0.09	± 0.57	± 0.79
CAVIA	57.25	66.25	67.50	68.50	36.72	6.01	3.99	3.27	1.99	1.27	1.21	1.17
	± 12.09	± 14.73	± 15.67	± 16.42	± 8.83	± 0.88	± 0.50	± 0.24	± 0.11	± 0.01	± 0.00	± 0.00
TAM (Comp)	63.00	76.50	82.75	88.5	6.15	3.43	2.69	2.13	2.33	1.30	1.23	1.19
	± 5.35	± 4.65	± 3.86	± 2.65	± 0.91	± 0.05	± 0.04	± 0.02	± 0.18	± 0.02	± 0.01	± 0.01
TAM (Non-comp)	45.25	72.5	81.5	89.75	7.80	5.08	3.60	2.42	4.37	1.27	1.17	1.11
	± 3.59	± 3.70	± 2.89	± 0.96	± 0.09	± 0.24	± 0.15	± 0.08	± 3.59	± 0.01	± 0.00	± 0.00

Table 8: Compositional models for few-shot sequence classification and sequence transduction. All models (except non-compositional TAM) get information on the primitives present in the tasks via extra tokens appended to the input sequence, except that one such primitive is unseen at test time. Non-compositional TAM is not given information about primitives, and estimates a single task embedding instead.