

**Supplementary Information for  
Quantized Memory-Augmented Neural Networks**

**A Fixed-Point Quantization and Computational Energy**

**A.1 Fixed-point quantization**

$U, V$  : floating-point vector

$\hat{U}, \hat{V}$  : fixed-point vector

$\hat{u}_i, \hat{v}_i$  : fixed-point scalar (binary vector)  $\in \{0, 1\}^n$

$\hat{u}_{ik}, \hat{v}_{ik} \in \{0, 1\}$

$\epsilon_{\hat{u}_i}$  : quantization error of  $\hat{u}_i$

$$\begin{aligned} \hat{u}_i &= u_i + \epsilon_{\hat{u}_i} \\ &= \text{sign} 2^{-FRAC} \sum_{k=0}^{n-2} 2^k \hat{u}_{ik} \end{aligned} \tag{8}$$

$$|\epsilon_{u_i}| < \begin{cases} 2^{-FRAC} & \text{if } |u_i| < 2^{IWL} \\ |2^{IWL} - |u_i|| & \text{if } |u_i| \geq 2^{IWL} \text{ (fixed-point overflow occurs)} \end{cases} \tag{9}$$

**A.2 Computational energy**

Table 3: Computation-energy gain from (Horowitz 2014)

Type	Arithmetic operation	Bit	Energy (pJ)	Gain <sup>a</sup>
Fixed point	add	8	0.03	123.3×
		32	0.1	37×
	mult	8	0.2	18.5×
		32	3.1	1.2×
Floating point	add	16	0.4	9.3×
		32	0.9	4.1×
	mult	16	1.1	3.4×
		32	3.7	1.0×

<sup>a</sup> compared with 32-bit floating-point mult

**B Experimental Results**

Table 4: Test error rates (%) of repeating each training 10 times for 8-bit fixed-point quantization on the bAbI dataset (COS=COSine similarity, HAM=HAMming similarity, ES=Early Stopping, MQ=Memory controller Quantization control)

Task	MANN												Q-MANN										
	floating point			fixed point (Q5.2)			fixed point (Q5.2)			fixed point (Q5.2)			fixed point (Q2.5)			fixed point (Q2.5)							
	COS	min	mean	std	COS	min	mean	std	COS	min	mean	std	ES	ES+MQ	COS	min	mean	std	HAM	ES	HAM	ES+MQ	
numerical representation																							
similarity measure function																							
1 1 supporting fact	0.0	0.00	0.000	0.0	16.25	34.259	0.0	2.09	5.274	0.0	1.19	0.936	0.1	0.97	0.801	0.3	0.73	0.483	0.4	1.01	0.465		
2 2 supporting facts	0.2	0.55	0.217	79.6	81.46	1.494	75.1	80	2.395	44.8	66.65	7.917	20.3	51.33	22.004	22.8	58.49	22.471	20.9	42.94	17.837		
3 3 supporting facts	23.1	25.47	1.282	78.5	80.71	1.942	78.5	82.11	3.198	72.3	76.58	2.686	75.7	79.77	2.739	74.8	77.5	1.671	53.7	68.28	5.923		
4 2 argument relations	30.8	32.08	1.049	31.6	33.03	1.370	31.3	32.67	0.730	40.7	48.83	5.997	31.4	32.77	1.656	30.5	33.23	3.774	30.6	35.8	6.806		
5 3 argument relations	13.5	15.08	0.750	69.6	76.6	6.130	68	70.14	0.969	39.8	46.82	7.262	13	14.66	1.152	12.6	15.26	1.581	15.3	16.44	0.635		
6 yes/no questions	3.1	5.23	2.591	49.3	49.94	0.389	47.9	49.57	0.617	47.7	49.58	0.807	15.1	17.81	2.433	15.6	20.64	4.160	12.1	15.56	2.334		
7 counting	11.4	12.63	1.041	51.1	52.83	0.790	22.1	30.97	8.749	21.7	23.04	1.227	21.3	24.57	3.192	20.8	21.8	0.872	17.8	19.3	1.166		
8 lists/sets	0.6	1.34	0.406	36.7	73.92	13.713	53.7	59.18	6.187	12	22.95	9.002	7.3	10.9	2.390	6.7	9.75	1.967	9.2	10.56	0.956		
9 simple negation	3.9	6.42	2.280	35.1	36.47	0.796	36.1	36.28	0.123	35.8	36.16	0.126	29.7	35.21	2.096	28.5	35.27	2.468	14.4	16.83	2.502		
10 indefinite knowledge	7.2	10.42	2.206	57.4	66.12	9.353	55.9	57.07	1.214	52.9	56.04	1.617	36.4	48.55	6.415	45.8	52.04	4.006	26.4	28.9	1.648		
11 basic coreference	0.2	9.11	5.572	11	31.82	25.406	9.6	10.85	0.712	9	11.81	1.735	1.9	10.72	3.485	1.7	13.66	5.571	7.5	11.69	2.788		
12 conjunction	0.0	0.00	0.000	0.0	16.82	35.435	0.0	0.43	1.290	0.0	8.25	10.995	0.0	0.29	0.260	0.0	5.7	17.007	0.0	0.6	0.467		
13 compound coreference	0.0	0.16	0.207	8.7	59.99	24.820	5.6	37.19	31.030	5.6	5.74	0.165	5.7	7.86	2.450	19.3	33.26	11.866	0.0	5.87	2.882		
14 time reasoning	3.3	3.90	0.313	5.5	20.18	21.807	4.7	20	16.633	9.8	17.78	3.331	14.1	15.56	1.083	13.1	15.13	1.084	13.9	28.33	13.601		
15 basic deduction	10.7	13.48	1.850	51.4	54.78	2.193	49.9	53.21	2.275	51	54.51	2.315	14.8	28.66	15.171	15.4	37.69	12.936	15.8	33.52	10.069		
16 basic induction	50.8	53.19	1.537	51.5	58.27	4.649	53	56	3.497	49.7	53.01	2.175	50.7	52.6	1.341	50.7	51.77	0.657	49.6	51.95	1.372		
17 positional reasoning	45	46.56	1.246	47.8	48.78	1.698	52	52	0.000	48	51.2	1.687	37.6	39.71	2.591	37.7	38.95	0.826	37.5	38.88	1.107		
18 size reasoning	38.4	42.32	2.185	46.2	48.58	2.553	52.9	53.57	0.564	42.4	47.43	3.952	41.5	45.85	2.570	42.9	45.96	1.893	40	43.21	1.596		
19 path finding	64.3	64.86	0.350	89.8	90.82	0.924	89.8	90.63	0.523	90.1	91.24	0.554	90	90.89	0.702	90.2	90.78	0.379	80.8	83.93	2.074		
20 agent's motivation	0.0	0.00	0.000	0.0	27.26	36.234	0.0	3.4	3.723	0.0	1.16	1.083	0.0	0.0	0.000	0.0	0.0	0.000	0.0	0.0	0.000		
Average error (%)	15.325	17.14	1.254	40.040	51.232	11.298	39.305	43.868	4.485	33.665	38.499	3.278	25.330	30.434	3.727	26.470	32.881	4.784	22.295	27.680	3.811		

Table 5: Test error rates (%) of repeating each training 10 times for 8-bit fixed-point and binary quantization on the bAbI dataset (COS=COsine similarity, HAM=HAMming similarity, ES=Early Stopping, MQ=Memory controller Quantization control)

Task	MANN												Q-MANN											
	floating point			fixed point (Q5.2)			fixed point (Q5.2)			fixed point (Q5.2)			fixed point (Q2.5)			fixed point (Q2.5)								
	COS	min	mean	std	COS	min	mean	std	COS	min	mean	std	ES	ES+MQ	COS	min	mean	std	HAM	ES	ES+MQ	HAM	min	mean
1 1 supporting fact	0.0	0.00	0.000	2.5	58.27	25.378	0.0	0.18	0.193	0.0	2.9	1.058	0.3	2.42	4.143	1.1	1.34	0.178	1.3	1.34	0.178	1.3	2.7	1.738
2 2 supporting facts	0.2	0.55	0.217	61.5	69.4	4.046	65.8	72.21	5.024	19.3	74.7	17.069	34.5	55.92	12.417	38.6	60.93	11.992	27.9	60.93	11.992	27.9	49.54	16.441
3 3 supporting facts	23.1	25.47	1.282	79	81.68	2.125	78.5	83.12	2.603	73.3	77.9	1.335	68.8	71.1	1.409	68.7	70.57	1.008	65.5	70.57	1.008	65.5	71.23	2.318
4 2 argument relations	30.8	32.08	1.049	34.6	45.55	10.995	35.6	39.86	2.449	35.3	61.2	7.573	42.2	47.37	3.144	44.5	49.16	3.356	41	49.16	3.356	41	44.84	2.056
5 3 argument relations	13.5	15.08	0.750	32.6	60.05	15.753	28.4	43.9	13.898	34.5	57.6	8.849	19.1	23.19	4.021	16.6	19.27	2.400	17.1	19.27	2.400	17.1	19.04	1.065
6 yes/no questions	3.1	5.23	2.591	22.9	40.45	10.501	24.5	37.89	10.597	48.1	52	1.079	14.9	18.7	4.225	13.4	19.58	8.248	14.2	19.58	8.248	14.2	22.8	6.186
7 counting	11.4	12.63	1.041	22.3	29.95	5.713	22	24.54	2.448	18	26.5	2.507	17.6	19.07	1.253	17.7	19.23	1.352	18.3	19.23	1.352	18.3	20.01	1.031
8 lists/sets	0.6	1.34	0.406	20.8	37.61	12.054	16.4	34.42	8.679	9.4	33.7	8.850	10.7	14.79	7.396	10.3	18.07	9.433	10.6	18.07	9.433	10.6	11.62	0.939
9 simple negation	3.9	6.42	2.280	21.6	28.16	4.146	26.3	29.3	3.765	23.3	28.6	1.611	13.2	16.22	1.550	12.2	17.33	2.261	14.9	17.33	2.261	14.9	16.89	1.610
10 indefinite knowledge	7.2	10.42	2.206	41.9	47.04	4.743	44.8	46.29	2.355	42.1	50.7	2.149	28.2	37.8	5.173	31.4	37.05	3.623	27.4	37.05	3.623	27.4	32	2.978
11 basic coreference	0.2	9.11	5.572	11.1	26.33	22.571	10.3	11.85	1.417	9.3	34.5	7.502	9.6	12.11	1.189	11.1	13.64	2.926	10.9	13.64	2.926	10.9	12.23	1.163
12 conjunction	0.0	0.00	0.000	0.0	6.93	14.217	0.0	0.06	0.084	0.0	3.2	1.043	1.1	3.49	1.994	1.4	16.28	15.443	1.2	16.28	15.443	1.2	7.7	10.675
13 compound coreference	0.0	0.16	0.207	5.6	47.08	29.219	5.6	5.6	0.000	5.4	8.6	1.165	4.9	6.45	1.439	4	6.55	1.912	0.2	6.55	1.912	0.2	5.51	2.204
14 time reasoning	3.3	3.90	0.313	10.2	14.85	4.348	8	14.93	6.156	10.6	19	2.615	31.5	40.02	4.653	28.3	39.04	7.043	27.2	39.04	7.043	27.2	35.6	7.773
15 basic deduction	10.7	13.48	1.850	23	52.78	10.862	30.5	54.18	8.635	47.6	57.3	2.629	29.8	44.67	5.847	19.8	39.56	11.884	17.5	39.56	11.884	17.5	42.48	9.933
16 basic induction	50.8	53.19	1.537	53	56.29	2.725	54.8	55.91	1.093	51.4	58.4	2.299	50.7	52.23	1.480	50.7	52.17	1.081	49.7	52.17	1.081	49.7	52.06	1.265
17 positional reasoning	45	46.56	1.246	44.5	47.33	2.110	47.8	49.33	1.140	49	52.6	1.036	37.4	39.64	1.796	36.6	40.09	2.938	38.6	40.09	2.938	38.6	40.62	1.971
18 size reasoning	38.4	42.32	2.185	44.7	47.24	2.689	44.8	50.01	3.609	42.8	51.1	2.482	41.3	45.19	1.846	42.7	45.95	1.317	44.1	45.95	1.317	44.1	45.89	1.658
19 path finding	64.3	64.86	0.350	84.4	87.42	2.410	84.1	86.94	2.482	83.4	90.1	1.810	86.4	87.54	0.929	84.1	87.25	1.325	84.9	87.25	1.325	84.9	87.15	1.023
20 agent's motivation	0.0	0.00	0.000	0.0	5.34	6.901	0.0	1.81	1.783	0.0	4.1	1.335	0.0	0.18	0.210	0.0	0.04	0.126	0.0	0.04	0.126	0.0	0.03	0.067
Average error (%)	15.325	17.14	1.254	30.810	44.488	9.675	31.410	37.117	3.921	30.140	42.235	3.800	27.110	31.905	3.306	26.660	32.655	4.492	25.625	32.655	4.492	25.625	30.997	3.705

## C MANN Model Description

Table 6: Model Descriptions

Symbol	Description	Domain
$I$	dimension of input	$\mathbb{N}$
$E$	dimension of internal representation	$\mathbb{N}$
$L$	number of memory element	$\mathbb{N}$
$R$	number of read	$\mathbb{N}$
$V$	input vectors (sentences)	$\mathbb{R}^{I \times L}$
$q$	input vector (question)	$\mathbb{R}^I$
$W_a$	weight of input(V)	$\mathbb{R}^{E \times I}$
$W_q$	weight of input(q)	$\mathbb{R}^{E \times I}$
$W_r$	weight of read memory	$\mathbb{R}^{E \times I}$
$W_k$	weight of read key	$\mathbb{R}^{E \times E}$
$W_o$	weight of output	$\mathbb{R}^{I \times E}$
$M_a$	address memory	$\mathbb{R}^{E \times L}$
$M_r$	read memory	$\mathbb{R}^{E \times L}$
$k_i$	$i$ th read key ( $1 \leq i \leq R$ )	$\mathbb{R}^E$
$w_{r,i}$	$i$ th read weight	$\mathbb{R}^L$
$r_i$	$i$ th read vector	$\mathbb{R}^E$
$o_i$	$i$ th output vector	$\mathbb{R}^I$

Memory addressing (content-based):

$$S(u, v) = u \cdot v$$

$$C(M, k)[i] = \frac{\exp\{S(M_i, k)\}}{\sum_j \exp\{S(M_j, k)\}}$$

Memory update:

$$M_a = W_a V$$

$$M_r = W_r V$$

Memory read:

$$k^i = \begin{cases} W_q q & \text{if } i = 1 \\ W_k k_{i-1} + r_i & \text{otherwise} \end{cases}$$

$$w_{r,i} = C(M_a, k_i)$$

$$r_i = M_r w_{r,i}$$

Output:

$$o_i = \text{softmax}(W_o k_i)$$

## D Analysis of the Effect of Quantization Error on Conventional MANN

### D.1 Vector similarity measure function - dot product

$$\begin{aligned} \hat{Z} &= \hat{U} \cdot \hat{V} \\ &= \sum \hat{u}_i \hat{v}_i \\ &= \sum (u_i + \epsilon_{u_i})(v_i + \epsilon_{v_i}) \\ &= \sum u_i v_i + \sum (u_i \epsilon_{v_i} + v_i \epsilon_{u_i}) + \sum \epsilon_{u_i} \epsilon_{v_i} \\ &\approx \sum u_i v_i + \sum (u_i \epsilon_{v_i} + v_i \epsilon_{u_i}) \\ &= Z + \epsilon_Z \end{aligned} \tag{10}$$

## D.2 Normalization function - Softmax

$$\begin{aligned}\hat{y}_i &= \frac{\exp(\hat{z}_i)}{\sum \exp(\hat{z}_k)} \\ &= \frac{\exp(z_i + \epsilon_{z_i})}{\sum \exp(z_k + \epsilon_{z_k})} \\ &= \frac{\exp(z_i)}{\sum \exp(z_k + \epsilon_{z_k} - \epsilon_{z_i})} \\ &\leq \frac{\exp(z_i)}{\sum \exp(z_k - \epsilon_{max})} \\ &= \frac{\exp(z_i)}{\exp(-\epsilon_{max}) \sum \exp(z_k)} \\ &= \exp(\epsilon_{max}) y_i\end{aligned}\tag{11}$$

## E Hyperparameters

Table 7: Hyperparameters

Parameter	Value
dimension of input ( $I$ )	17 - 98 <sup>a</sup>
dimension of internal representation ( $E$ )	60
number of memory locations ( $L$ )	50
number of read ( $R$ )	3
learning rate	0.3
Hamming similarity weight constant ( $\alpha$ )	-3

<sup>a</sup> depend on the task in the bAbI dataset