

1 We thank the reviewers (**R4, R5, R6**) for the useful feedback. Below we address their questions.

2 **(R4, R5, R6) Is STARCAPS a capsule network?** The table to the  
3 right shows the test accuracies of two types of experiments on Small-  
4 NORB (novel, familiar viewpoints). **Type1:** 3 runs of models EMCaps  
5 {64, 8, 16, 16, 5}, STARCAPS {32, 8, 8, 8, 5}, fully trained on familiar  
6 views and tested on both novel and familiar views. **Type2:** EMCaps  
7 {32, 32, 32, 32, 5}, STARCAPS {32, 32, 16, 16, 5}, trained on familiar

Model #params	Type1 (low capacity models)		Type2 (high capacity models)		
	EMCaps 68K	STARCAPS 73K	CNN 4.2M	EMCaps 316K	STARCAPS 318K
Familiar	95.66±0.03	95.72±0.02	96.3	96.3	96.3
Novel	86.12±0.05	86.07±0.03	80.0	86.5	86.3

8 views and early stopped when test accuracy reached 96.3% (as the CNN model in Table 2 of [EMCaps, Hinton et al]).  
9 In *Type1*, we notice that STARCAPS achieves comparable results (small difference in accuracy) to EMCaps both on  
10 familiar and novel viewpoints. In *Type2*, on the novel viewpoints, STARCAPS performs dramatically better than CNN  
11 model (+6.3%) and its accuracy is only slightly lower than EMCaps (-0.2%).

12 affNIST results: We trained STARCAPS {32, 8, 16, 16, 10} on MNIST following the data augmentation mentioned by  
13 the authors on OpenReview/ICLR18. The test accuracy of STARCAPS on affNIST is 93.03% vs. 93.1% for EMCaps  
14 {32, 32, 32, 32, 10}. In conclusion, the results show that STARCAPS is capable of detecting novel viewpoints similarly  
15 to EMCaps, retaining capsules properties.

16 **(R4) Performance of STARCAPS vs. CNNs** Although the main purpose of STARCAPS is to alleviate the computa-  
17 tional complexity of baseline capsule networks while being able to detect viewpoint variations, STARCAPS models  
18 achieve accuracies nearly on par with those modern CNN models. On CIFAR10, STARCAPS: 91.23%, #params=80K  
19 vs. ResNet20: 91.25%, #params=270K vs. ResNet110: 93.57%, #params=1.7M. On CIFAR100, STARCAPS: 67.66%  
20 vs. ResNet38: 68.54% vs. ResNet110: 71.21%. It is possible that scaling up STARCAPS models to match #params in  
21 ResNet, would lead to better performance; however this requires further extensive study.

22 **Ablation studies** (a) *ST-Router*: Removing ST-Router leads to lower performance. On MNIST, STARCAPS model  
23 {32, 8, 16, 16, 10} achieves 99.41 w/o ST-Router & 99.59 with ST-Router, while {32, 4, 64, 4, 10} achieves 98.37  
24 w/o ST-Router & 99.48 with ST-Router. (b) *Single parent assumption*: The single parent assumption enforced in  
25 EMCaps/DynamicCaps, while it may allow better encoding of entity representation, it imposes a limitation as it ignores  
26 actual/natural use cases for object recognition. We tested STARCAPS models designed to force the single parent assump-  
27 tion, the results were comparable to the proposed STARCAPS models on MNIST; however on {CIFAR10; CIFAR100}  
28 the results were inferior due to varied clutter in backgrounds, {89.91; 62.33} vs. STARCAPS = {91.23; 67.66} &  
29 EMCaps = {88.10;  $n/a$ }. (c) *Effect of sharing weights & role of attention*: Experiments with two settings. First, STAR-  
30 CAPS with separate weights with Attentions. We didn't notice improvement on MNIST; on CIFAR10 {32, 8, 8, 8, 10}  
31 achieved 91.31 vs. 91.23. However, the train/test time were significantly higher due to extra matrix multiplications  
32 as in EMCaps; we couldn't train models on CIFAR100. Second, STARCAPS with separate weights w/o Attentions.  
33 Experiments on MNIST/CIFAR10 showed very poor performance. In conclusion, the proposed setting of STARCAPS  
34 provides best results in general, in terms of accuracy and train/test time while preserving capsule properties.

35 **(R5) Visualization of instantiations params** We will include visualizations in the supplementary material.

36 **(R5) Vanishing pose** In EMCaps the initializations of vector parameters ( $\beta_a, \beta_u$ ) control the initial sparsity. In practice,  
37 according to our experiments, even with careful initialization of parameters, some EMCaps models suffer from unstable  
38 performance due to numerical issues with gradients (vanishing gradients). This was confirmed by the authors of EMCaps  
39 (see OpenReview/ICLR18 comments). We noticed unstable performance in EMCaps when a capsule layer has very  
40 large #capsules compared to lower/higher capsules, and when multiple adjacent layers have very large #capsules. The  
41 routing in STARCAPS automatically prunes the unneeded capsules without being sensitive to #capsules/initializations.  
42 We will update Table 1, fix the argument in lines (264-280), and add clarifications about the instability issues.

43 **(R6) SmallNORB overall results** STARCAPS: **S1**= {32, 8, 8, 8, 5}, 73K, 98.0%; **S2**= {32, 32, 16, 16, 5}, 318K,  
44 98.2%; vs. EMCaps: **E1**= {64, 8, 16, 16, 5}, 68K, 97.8%, **E2**= {32, 32, 32, 32, 5}, 316K, 98.2%

45 **(R6) Related work, capsules introduction, pseudo code** We will add pseudo code in the supplementary material,  
46 references to (Zhang et al.) and (Yang et al.), and a more detailed introduction to capsule networks.

47 **(R6) Global avg pooling (GAP)** In STARCAPS, the role of GAP is not routing. We use GAP in Decision-Learner  
48 internally in ST-Router ( $\mathcal{R}_{ij}$ ) to rapidly capture confidence maps from an attention matrix  $A_{ij}$ . The role of  $\mathcal{R}_{ij}$  is to  
49 estimate binary connectivity decision signal between two capsules; each  $\mathcal{R}_{ij}$  determines the connectivity between a  
50 single lower-level capsule and a single higher-level capsule (one-to-one), and not routing between lower and higher  
51 capsules (compared to dynamic routing in capsules, or static routing between neurons using max-pooling in CNNs).  
52 The name "ST-Router" may raise confusion with "routing" in EMCaps/DynamicCaps. Each ST-Router acts as a *gating*  
53 mechanism between two capsules, whereas the whole set of ST-Routers acts as a *routing* mechanism between all  
54 capsules. We will change the name to "ST-Gate". We also use GAP after the final output layer (CLASSCAPS) to calculate  
55 the final activation probabilities from the sigmoid of pose matrix, and not for routing between capsules.

56 **(R6) Hard attention** STARCAPS uses both soft and hard attention modules. The soft attention (Attn-Estimator)  
57 estimates soft relevance signal for each higher capsule which is used for scaling the pre-vote. To sparsify the network  
58 we use a hard attention module (ST-Router). Each route can be seen as a double-attention (soft & hard) mechanism.