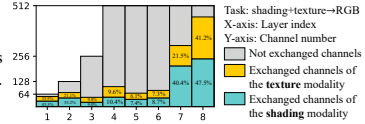


Table A: More ablation studies on dataset NYUDv2 with ResNet101. All experiments use unshared BNs. Ensemble indicates using ensemble to mix predictions of RGB&Depth.

Method	Convs	Ensemble	Mean IoU
Summing activation	Shared	✓	48.1 (%)
Summing activation + soft gate	Shared	✓	48.4 (%)
CBN	Unshared	✓	48.3 (%)
CBN	Shared	✓	48.9 (%)
Ours (concat-fusion block)	Shared	×	50.8 (%)
Ours	Shared	✓	51.1 (%)

Figure A: Percentage of exchanged channels on dataset Taskonomy. Task: shading+texture→RGB. X-axis: Layer index. Y-axis: Channel number. Legend: Not exchanged channels (grey), Exchanged channels of the texture modality (yellow), Exchanged channels of the shading modality (blue).



1 We thank all reviewers for their constructive comments and address the reviewers’ concerns point by point below.
2 **[To Reviewer #1:] Q1. More mid-fusion baselines:** Thanks for the suggestions. We have conducted extra ablation
3 studies raised by the reviewer in Table A. For the CBN strategy [DeVries et al. 2017] specifically, we modulate the BN
4 of one modality conditional on the other. In this sense, CBN performs cross-modal message passing via BN modulation,
5 which is clearly different from our method that directly exchanges channels between different modalities for fusion. As
6 expected, all ablations (including CBN) are inferior to our method, which again verifies the validity of the proposed
7 mechanism. Per the reviewer’s concern, experiments on other datasets will also be added into the final version.
8 **Q2. More explanations on fundamental parts: 1.** We can **not** perform channel exchanging without ℓ_1 , as it enables
9 the discovery of unnecessary channels and comes as a pre-condition of Theorem 1. Naively exchanging, *e.g.*, 30%
10 channels only gets IoU 47.2. **2.** We provide the comparison with the modulation approach in Table A. **3.** We can conduct
11 channel exchanging upon unshared CNN, which has been evaluated in Table 6 (the 5th and 7th rows). Promisingly, its
12 full-channel case getting 49.1 still outperforms baselines in Table 2 under the same setting. We believe our method is
13 generally useful and channels are aligned to some extent given unshared CNN. **4.** We have discussed why regularizing
14 half of the channels is beneficial in L155-161 and provided its empirical evidence in Table 1 (the 5th and 6th rows).
15 **Q3. The rationality of our theorems: Theorem 1** is meaningful and crucial. Yes, ℓ_1 makes the parameters sparse, but
16 it can not tell if each sparse parameter will keep small in training considering the gradient in Eq.(4). Conditional on BN,
17 Theorem 1 proves that $\gamma = 0$ is attractive, which is nontrivial and novel. **Corollary 1** states that f' is more expressive
18 than f when $\gamma = 0$, and thus the optimal f' always outputs no higher loss, which, yet, is not true for arbitrary f' (*e.g.*
19 $f' = 10^6$). Besides, Corollary 1 holds upon unshared CNN (see L191) and is consistent with Table 6 in the unshared
20 scenario (full-channel: 49.1 vs half-channel: 48.5), although full-channel exchanging is worse under the sharing setting.
21 **Q4. On recommended references:** We thank for the comment, will cite all the raised references (including the
22 modulation-based and CBN methods), and will refresh the references related to early/late fusion and the fusion
23 taxonomy. Our claims of “Introduction 3rd paragraph” are intuitively supported by the observations from Figure 1,
24 namely, the aggregation-based fusion integrates multiple networks into one single branch and the alignment-based
25 fusion conducts an indirect way of fusion by training. We will further clarify this in the final version.
26 **Q5. Other comments: 1.** In L201, “the standard setting” refers to train/test split, while “commonly/same with our
27 setting” is on architecture design, as already described in L236-238. **2.** We agree our current version is inapplicable for
28 heterogeneous fusion (*e.g.* CNN+RNN). That is why we, explicitly and honestly, point it out in L170-176. Given the
29 effectiveness on homogeneous fusion, we conjecture that adding extra homogeneous MLP layers upon each CNN/RNN
30 enables channel (or neural unit) exchanging for heterogeneous fusion, which will be left for future exploration.
31 **[To Reviewer #2:] Q1. More experiments:** We have carried out random exchanging like ShuffleNet or directly
32 discarded unimportant channels without channel exchanging on NYUDv2, the IoUs of which are 46.8 and 47.5,
33 respectively. Both are much worse than our method (51.1). We are willing to add these ablations into the final version.
34 **Q2. Why use the mean:** We have tried weighted summation of channels, but do not observe clear improvement (FID
35 from 60.90 to 60.94 on the translation task Depth+Normal+Texture→RGB). We conjecture that γ of each modality in
36 Eq. (6) already contributes to the weights, thus the simple and parameter-free mean delivers promising performance.
37 **[To Reviewer #3:]** We thanks for the positive comments, and are sorry for the unclear presentations on Theorem 1.
38 We would like to highlight that our conclusion always holds regardless of the value of β . Our proof (in appendix) is
39 interested in the gradient $dL/d\gamma$ only around $\gamma = 0$ (not other points) and finds that the probability of staying around
40 $\gamma = 0$ is large and hits the Gaussian peak, even for the extreme case (β or dL/dx' is non-zero) raised by the reviewer.
41 By mentioning $\beta = 0$ in L179, we initially tend to simplify the analysis of Corollary 1. Yet, our current proof of
42 Corollary 1 does not require this assumption, and we will remove $\beta = 0$ to avoid confusion in the final version.
43 **[To Reviewer #4:] Q1. Novelty and contributions compared to [38]:** Our paper is inspired but differs from [38] in
44 two aspects: **1.** while [38] aims at model compression, our paper attempts to conduct multimodal fusion. Thus upon
45 [38], we employ channel exchanging for message passing, which is simple yet effective and novel compared to current
46 fusion methods. **2.** Theorem 1 is related to [38], but [38] never confirms that unnecessary channels during training can
47 hardly recover, and we rigorously prove it by exploring the behaviors of the ℓ_1 norm and BN layers.
48 **Q2. More discussions: 1.** Thanks for the reminding, and we have reported the percentage of exchanged channel for
49 every layer in Figure A. It shows that upper layers are more active in exchanging. **2.** Please refer to Q5.2, R#1 for the
50 discussion on heterogeneous fusion. **3.** Regarding the three-modal example (A, B, and C) raised by the reviewer, the
51 ℓ_1 norm is conducted on disjoint channels of different modalities. It means in Eq.(6) that for each channel, only one
52 modality will meet the criterion (scaling factor lower than the threshold). This setting will avoid unavailing exchanging
53 as already described in L160-161. **4.** Not all experiments specify two modalities. We have tested the effectiveness of
54 our method with modalities more than two in Table 5 (also Table 7, 8 in appendix), where our method still makes effect.
55 **Q3. Other comments:** We use V100 with GPU memory 32G, and the segmentation needs 20G. We will detail the
56 GPU type in the final version. Applying the ℓ_1 norm alone delivers no improvement, but it is necessary for the discovery
57 of unnecessary channels and channel exchanging (see also our response in Q2.1, R#1).