
Disentangled Counterfactual Learning for Physical Audiovisual Commonsense Reasoning *Supplementary Material*

Anonymous Author(s)

Affiliation

Address

email

1 The supplementary material provides detailed implementation information on the baselines used for
2 comparison, as well as various analyses, such as hyper-parameters, model size and training time,
3 audio disentanglement learning, static factors, dynamic factors, and physical knowledge relationships.
4 Moreover, we show more visualization results in experiments.

5 1 Implementation Details of Compared Baselines

6 **LateFusion.** Following PACS [1], we used a pre-trained Debert-a-V3-Large [2] as the text encoder
7 to encode all questions into \mathbb{R}^d , where $d = 768$, which was saved during the training. We did not
8 apply any data augmentation to the text and extracted the text embeddings from the <CLS> token of
9 the text model’s output layer (pre-pooler). The pre-trained model can be downloaded from here¹.

10 For videos, we downsampled them as input for the model. We used the ViT/B-16 model pre-trained
11 on ImageNet-21k provided by HuggingFace² to extract features from the video frames. Following
12 the video augmentation steps used in the pre-trained model, we began with video frames of size 252
13 $\times 252$ and randomly selected 8 evenly spaced frames. We then cropped the same 224×224 from
14 each frame and randomly flipped the images horizontally with a probability of 0.5.

15 For audio, we used the pre-trained AST (Audio Spectrogram Transformer) [3] model with a time and
16 frequency stride of 10 and weight averaging that was pre-trained on the full AudioSet [4]. The model
17 can be downloaded from here³. We followed the audio augmentation steps for pre-trained AST on
18 AudioSet [4] and PACS [1] using frequency and time masking [3]. We utilized 128 mel bins with a
19 target length of 1024. Then, we masked a band of size 48 in the frequency domain and a crew of size
20 144 in the time domain. Finally, we normalized the spectrogram as $\text{spec} = (\text{spec} + 4.26)/(4.57*2)$ and
21 added random noise.

22 During training, we use a simple grid search and set the learning rate to $5 \cdot 10^{-4}$, the weight decay to
23 $5 \cdot 10^{-5}$ and the batch size to 64 (20GB of GPU memory was used). We trained the model for 40
24 epochs with early stopping, and we freeze all backbone layers and only use trainable MLPs layers to
25 fuse multimodal information.

26 **CLIP.** For CLIP [5], we only used the video frames $X^v = \{X_1^v, X_2^v, \dots, X_T^v\}$ as input, where T
27 refers to the number of video frames and used the same ViT/B16 as the backbone network, along with
28 the same video preprocessing and augmentation, to obtain features $X = [X_s^v, X_z^v]$ for objects.
29 To ensure a fair comparison, we used the fusion and optimization method as same as Latefusion.
30 Additionally, we use learning rate of $1 \cdot 10^{-4}$ with weight decay of $1 \cdot 10^{-5}$ after a simple grid search.

¹<https://huggingface.co/microsoft/deberta-v3-large>

²<https://huggingface.co/google/vit-base-patch16-224-in21k>

³<https://github.com/YuanGongND/ast>

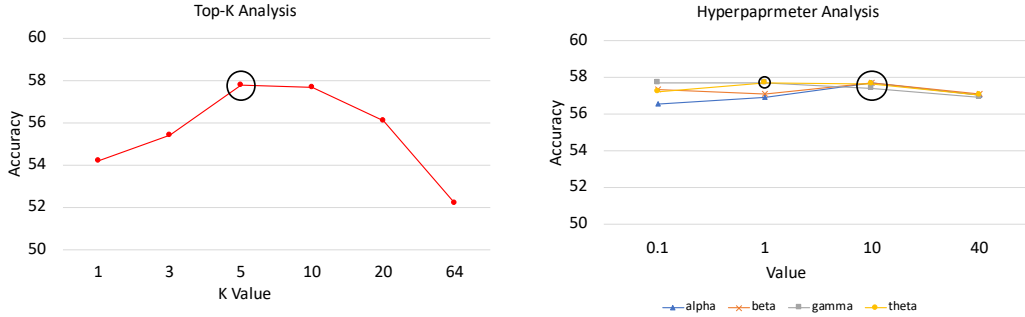


Figure 1: Hyper-parameters Analysis of our DCL.

31 **AudiCLIP.** For AudioCLIP [6], we used both audio X^a and video X^v as input. To ensure a fair
 32 comparison, we used the preprocessing and augmentation method as same as Latefusion. For
 33 hyperparameters, we used the same configuration as CLIP. We used ESRestNet(X)t-fbsp [7] as the
 34 backbone network for audio feature extraction. During training, we froze all layers of the backbone
 35 network, and only use trainable MLP layers to fuse multimodal information.

36 **UNITER.** In experiments, we utilized UNITER pre-trained on NLVR2 dataset [8] for feature
 37 extraction and fusion. The pre-processed sequential video information X^v was used as input. In
 38 our experiments, UNITER’s pair setup was applied to handle the input object-question pair as two
 39 independent text-video pairs. The concatenation of the output of both [CLS] in UNITER was regarded
 40 as \hat{Y}_{X,A_X} , as mentioned in the main paper. We use the learning rate of $1 \cdot 10^{-5}$ and a weight decay
 41 of 0.01.

42 2 Hyper-Parameters Analysis

43 Our proposed model consists of two main modules: Disentangled Sequence Encoder (DSE) and
 44 Counterfactual Learning Module (CLM). Specifically, in DSE, we used a hidden layer size of 256
 45 for Bi-LSTM and set $\gamma = 1$, $\alpha, \beta = 10$, and $\theta = 1$. In CLM, $\tau = 2$ and $k = 5$ were used when
 46 calculating similarities and constructing the physical knowledge relationships. The hyper-parameters
 47 for both DSE and CLM were kept consistent when incorporated with each baseline. Figure 1 shows
 48 the hyper-parameter analysis of Late Fusion w/ DCL, where the left part illustrates the relationship
 49 between the Top-K value and the corresponding accuracy, with k taking values of 1, 3, 5, 10, 20,
 50 and 64. When k=1, it means that the object’s physical properties are only related to itself, while
 51 k=64 represents the physical properties of the object that are related to all objects within the batch.
 52 From the figure, it can be clearly observed that when we set K as 5 the model can achieve the best
 53 performance, indicating that an appropriate K can help improve the ability to explore the common
 54 physical properties. Meanwhile, a value of k that is too large can introduce excessive noise while
 55 too small is insufficient to exploit the relevance among objects and both result in decreased accuracy.
 56 The right part in Figure 1 shows the results when different values of α , β , γ , and θ are adopted,
 57 demonstrating that the model is insensitive to these hyperparameters, indicating its robustness. The
 58 black circular markers in the figure indicate the parameter values that were ultimately used.

59 3 Model Size and Training Time

60 Table 1 reports the model size and training time of our proposed method and baselines, by using an
 61 Intel 6226R CPU and an NVIDIA RTX 3090 GPU, in terms of the time required for a single epoch of
 62 training. It can be observed that the increase in training time is acceptable, and the additional memory
 63 usage is also within a controllable range.

64 4 Analysis of Audio Disentanglement Learning

65 In this section, we will analyze the effectiveness of our proposed Disentangled Sequential En-
 66 coder (DSE) on the input audio data. As described in Section 3.1 in our paper, we represent audio

	Model Size	Train Time
Latefusion	170.2M	1,214s
Latefusion w/ DCL	189.4M	1,317s
AudioCLIP	230.1M	1,545s
AudioCLIP w/ DCL	242.6M	1,628s

Table 1: The analysis of the model size and training time

Baseline Model	Accuracy (%)	
	PACS	PACS-Material
Late Fusion [9]	55.0 ± 1.1	67.4 ± 1.5
Late Fusion w/ DSE-audio	56.9 ± 0.5	68.1 ± 0.4
AudioCLIP [6]	60.0 ± 0.9	75.9 ± 1.1
AudioCLIP w/ DSE-audio	61.5 ± 0.8	76.0 ± 0.7

Table 2: Performance comparison between our proposed DSE-audio and existing baseline methods.

67 features as a sequence and extract them as sequence data represented by $X^a = \{X_1^a, X_2^a, \dots, X_T^a\}$,
68 where T denotes the number of audio time steps. It should be noted that, in the audio decoupling
69 experiment, we did not decouple the video but processed it through averaging.

70 **Quantitative Results.** As shown in Table 2, we compare our method with other baseline methods.
71 Since only Latefusion [9] and AudioCLIP [5] take audio as input, we compared our method with
72 both of them. It can be observed from Table 2 that using only DSE-audio, Late Fusion achieved an
73 absolute improvement of 1.9%, while AudioCLIP achieved an absolute improvement of 1.5% on the
74 PACS dataset. This indicates that our proposed DSE-audio method has a significant impact on audio
75 decoupling. The same conclusion can also be drawn from the results on the PACS-Material dataset,
76 which demonstrates the effectiveness and superiority of our DSE as a plug-and-play method that can
77 achieve excellent performance on various datasets.

78 5 More Visualization Results

79 As shown in Figure 2 and Figure 3, we present more visualization results comparing our proposed
80 method with other baseline models. It can be seen from the figures that our proposed DCL method
81 outperforms the original method.

82 6 Analysis of Material and Question Properties

83 Figure 4 and Figure 5 show the results of object accuracy for per materials and per-property related
84 to each question, respectively. From Figure 4, we can see that after using the proposed DCL, there
85 was an improvement in object accuracy for all materials, especially for *Rubber*, *Plastic*, and *Metal*
86 show the most prominent improvement. This is because our proposed DCL is capable of capturing
87 clean dynamic features, which is essential in learning the physical properties of objects that often
88 require dynamic movements. Based on Figure 5, after incorporating our proposed DCL, the accuracy
89 of all properties related to questions improved, especially for Flexibility and Weight properties s
90 the most significant improvement. The reason may attribute to the fact that these questions require
91 more dynamic features for accurate judgment, and therefore, the usage of our proposed DCL led to a
92 significant improvement in performance.

93 7 Analysis of Static Factors

94 Figure 7 and Figure 8 show the results of our proposed method using only static factors and other
95 baselines. As shown in Figure 7, it can be observed that our model performs better than the baseline
96 in terms of “Stone”, “Glass”, and “Textiles” while underperforming on “Metals”, probably stems
97 from that metals require more audio information for support. This finding demonstrates the advantage
98 of the decoupled static factors in material classification. Figure 8 illustrates the performance of our

99 model on questions related to different physical properties, where it performs well w.r.t “texture”,
100 “shape”, and “size”, indicating the helpfulness of the decoupled static factors in these properties.

101 **8 Analysis of Dynamic Factors**

102 In Figure 6, we show a few additional examples of clustering using dynamic factors. It can be
103 observed that the decoupled dynamic factors represent similar or related action information.

104 **9 Analysis of Physical Knowledge Relationships.**

105 Figure 9 presents the visualization of the physical knowledge relationships captured by our dynamic
106 factors with the affinity matrix A , and the corresponding top-5 results are displayed. As shown in the
107 figure, we can find that the selected top-5 video in the top row and the bottom row in the red box
108 are actions that are similar or related to the given sample video, indicating the effective discovery of
109 common physical relevance through our physical knowledge mining, as mentioned in our main paper.

110 **References**

- 111 [1] Samuel Yu, Peter Wu, Paul Pu Liang, Ruslan Salakhutdinov, and Louis-Philippe Morency. Pacs:
112 A dataset for physical audiovisual commonsense reasoning. In *Computer Vision–ECCV 2022:
113 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVII*,
114 pages 292–309. Springer, 2022.
- 115 [2] Pengcheng He, Jianfeng Gao, and Weizhu Chen. Debertav3: Improving deberta using electra-style
116 pre-training with gradient-disentangled embedding sharing. *arXiv preprint arXiv:2111.09543*,
117 2021.
- 118 [3] Yuan Gong, Yu-An Chung, and James Glass. Ast: Audio spectrogram transformer. *arXiv preprint
119 arXiv:2104.01778*, 2021.
- 120 [4] Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing
121 Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for
122 audio events. In *2017 IEEE international conference on acoustics, speech and signal processing
123 (ICASSP)*, pages 776–780. IEEE, 2017.
- 124 [5] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
125 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
126 models from natural language supervision. In *International conference on machine learning*,
127 pages 8748–8763. PMLR, 2021.
- 128 [6] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip
129 to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics,
130 Speech and Signal Processing (ICASSP)*, pages 976–980. IEEE, 2022.
- 131 [7] Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Esresne (x) t-fbsp: Learning
132 robust time-frequency transformation of audio. In *2021 International Joint Conference on Neural
133 Networks (IJCNN)*, pages 1–8. IEEE, 2021.
- 134 [8] Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. A corpus of natural language for visual
135 reasoning. In *Annual Meeting of the Association for Computational Linguistics*, 2017. URL
136 <https://api.semanticscholar.org/CorpusID:19435386>.
- 137 [9] Yagya Raj Pandeya and Joonwhoan Lee. Deep learning-based late fusion of multimodal in-
138 formation for emotion classification of music video. *Multimedia Tools and Applications*, 80:
139 2887–2905, 2021.

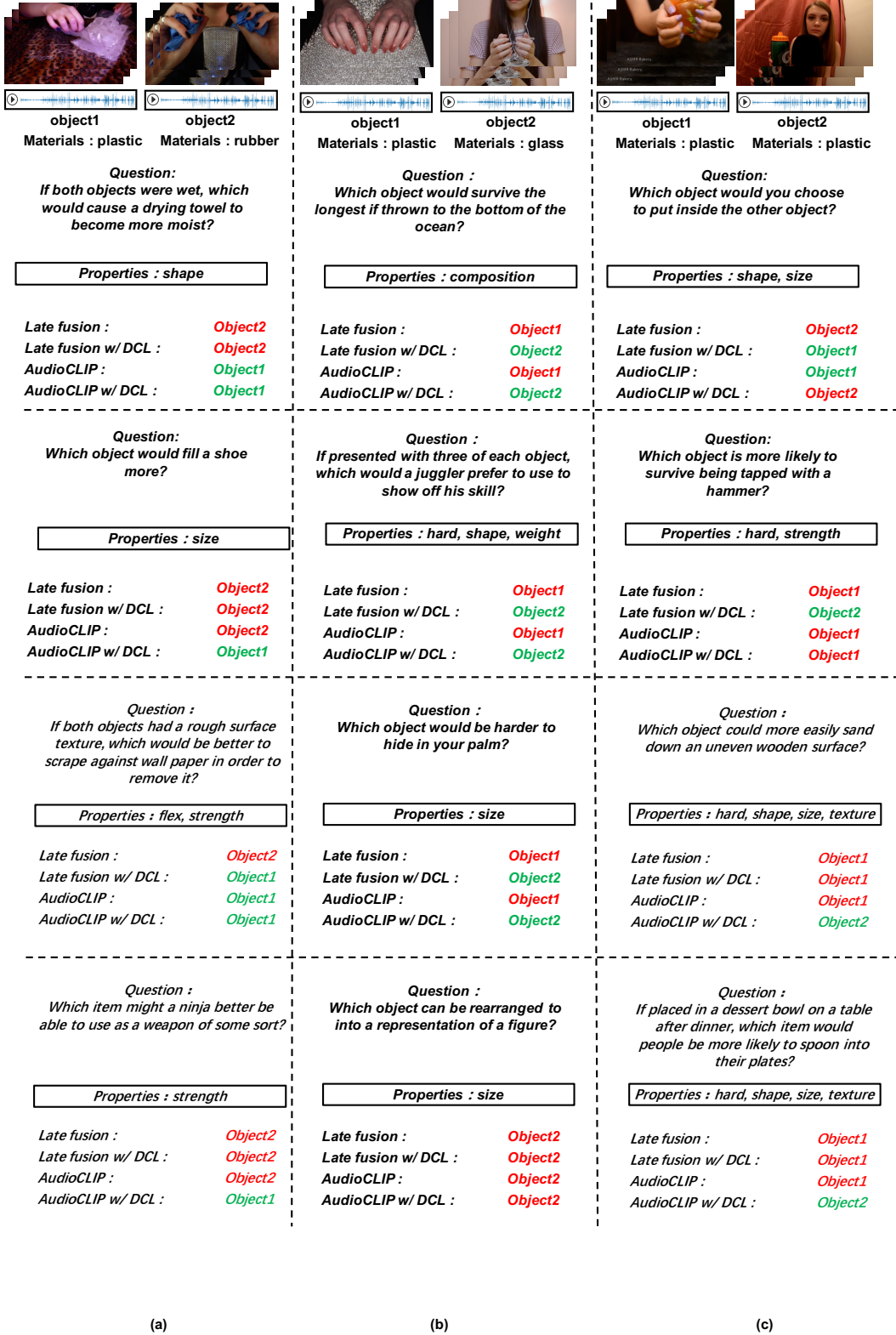


Figure 2: Comparison between our proposed method and existing baseline methods, where ‘w/ DCL’ indicates the baseline incorporated with our proposed DCL method.

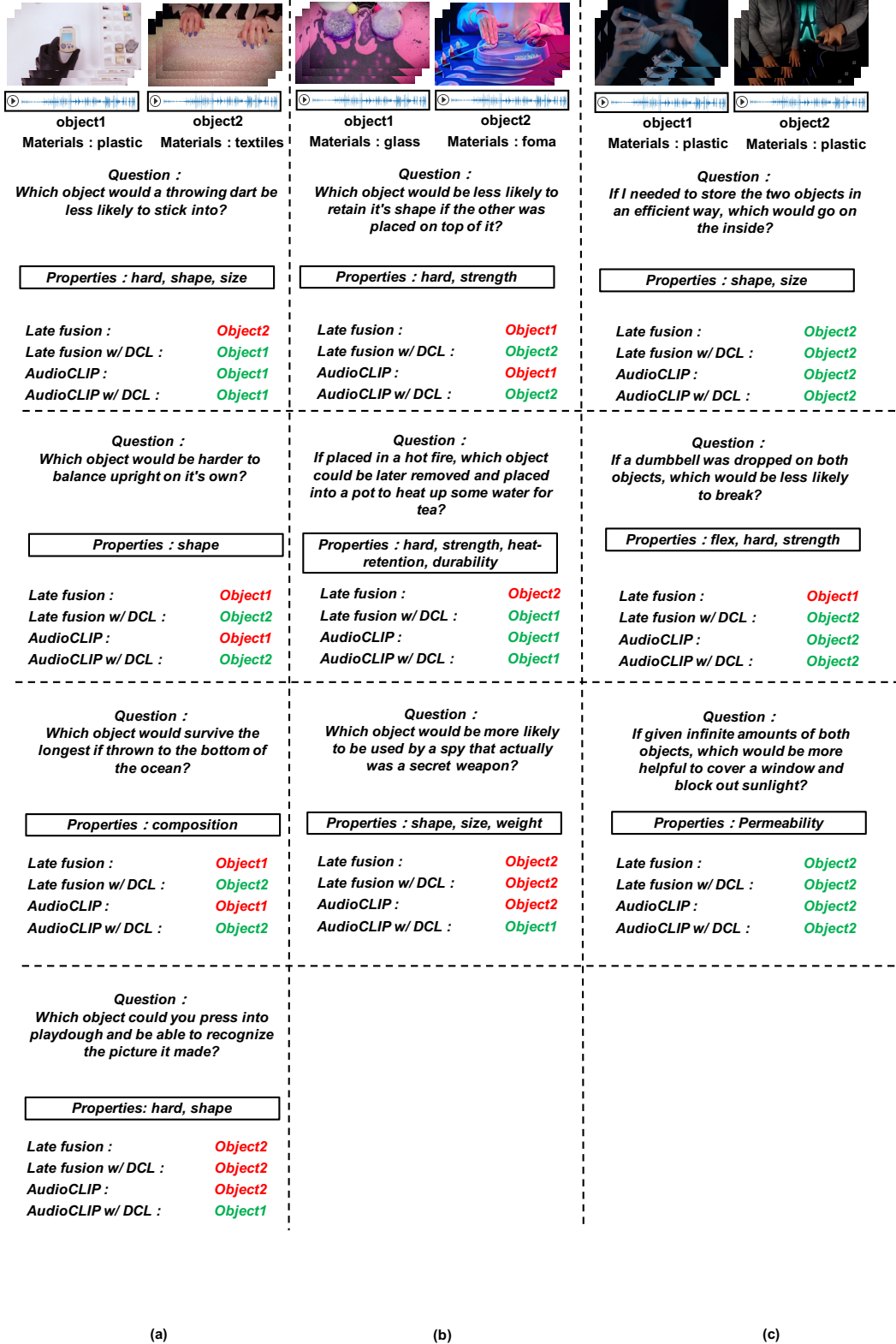


Figure 3: Comparison between our proposed method and existing baseline methods, where ‘w/ DCL’ indicates the baseline incorporated with our proposed DCL method.

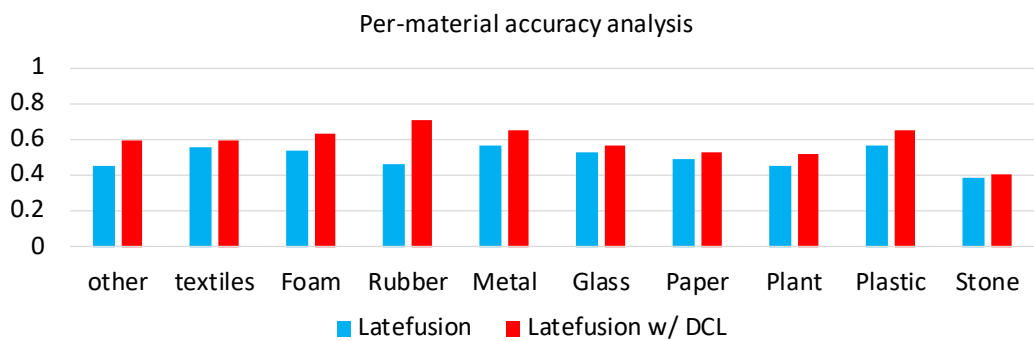


Figure 4: Accuracy results of per material object.

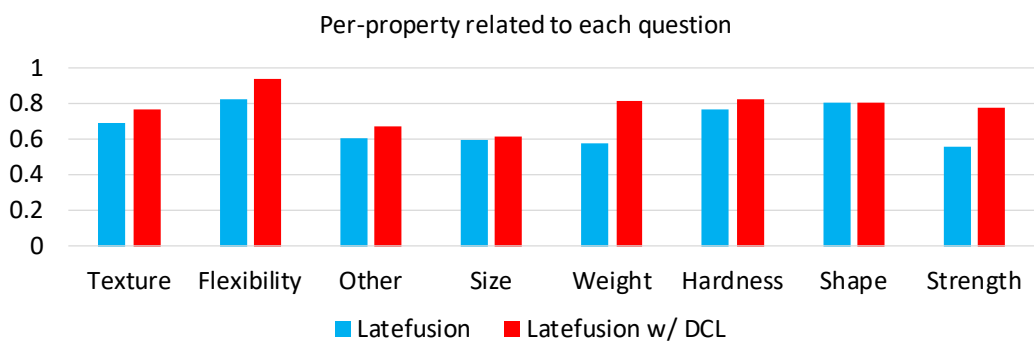


Figure 5: Accuracy results of per properties related to each question.

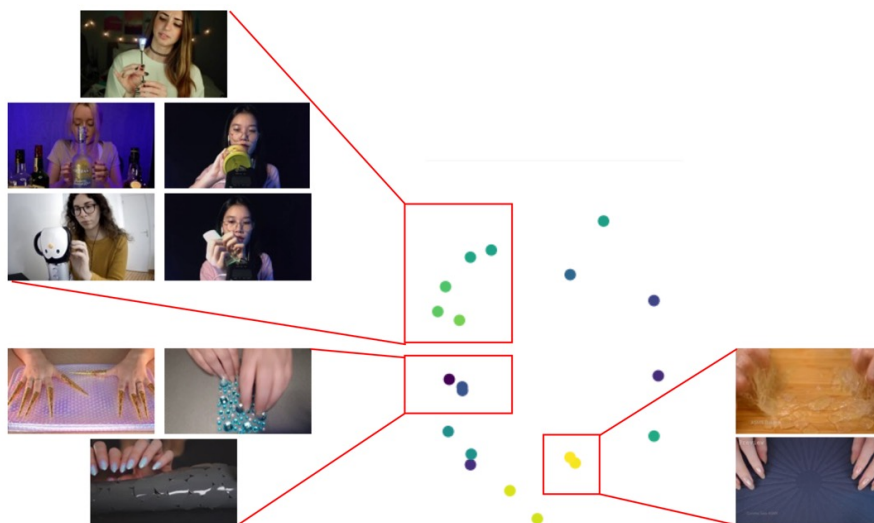


Figure 6: The t-SNE visualization of the obtained dynamic factors in the latent space of video samples in the test set.

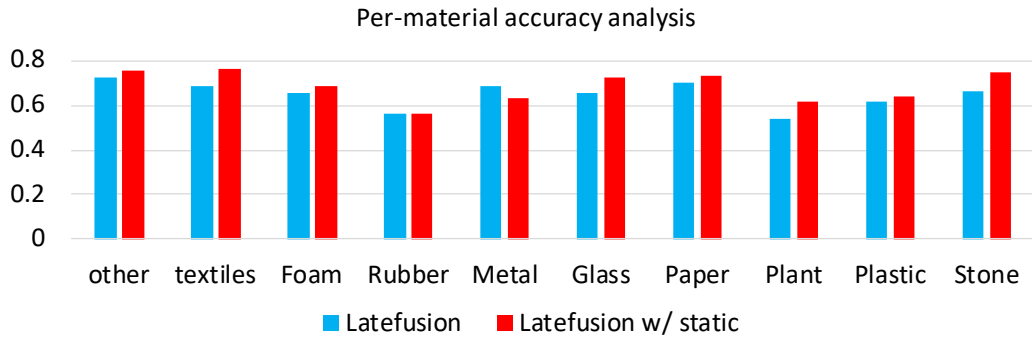


Figure 7: Accuracy results of per material object.

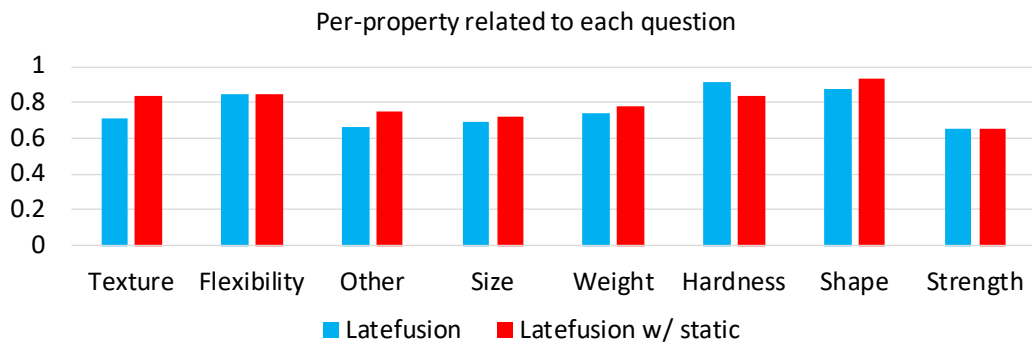


Figure 8: Accuracy results of per properties related to each question.

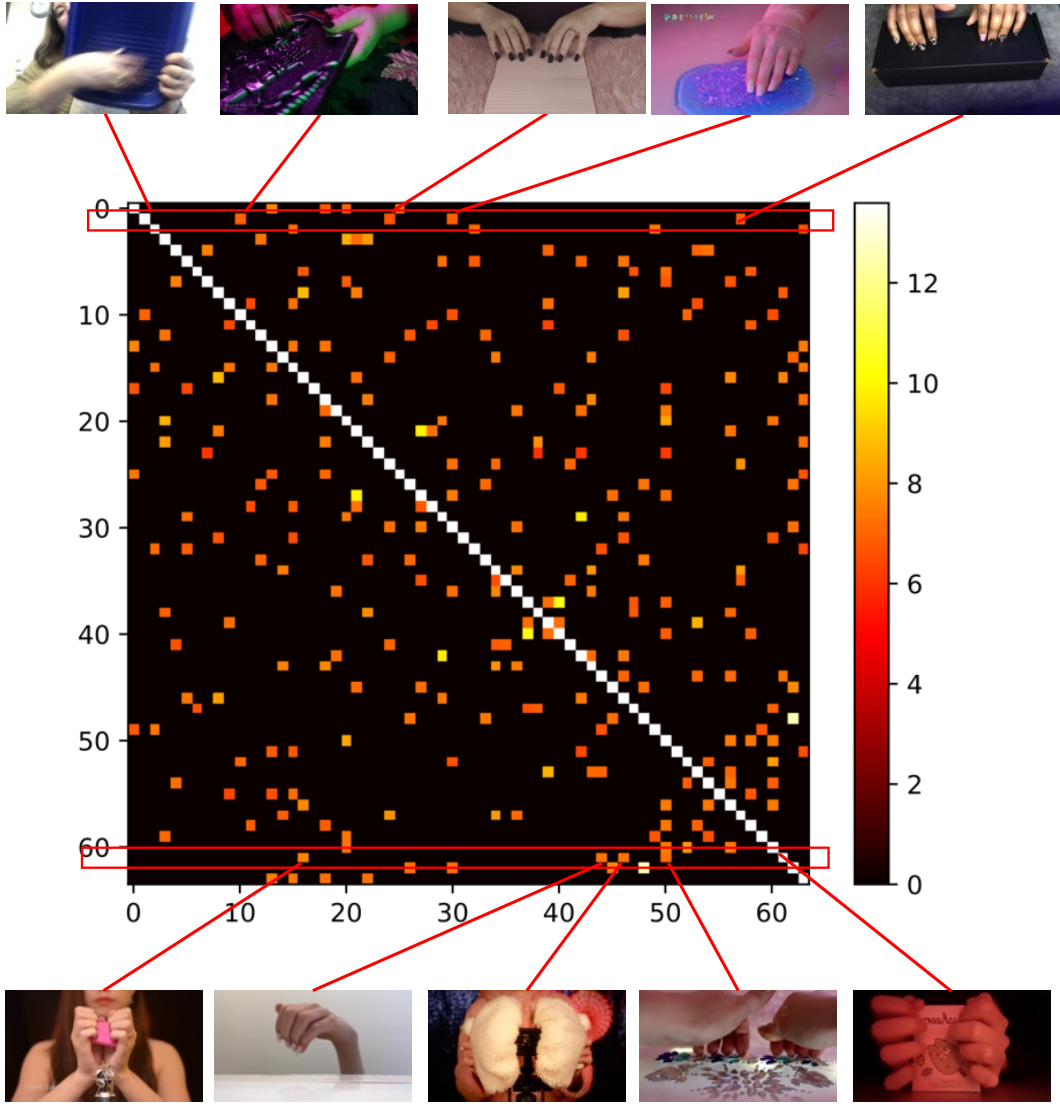


Figure 9: Visualization of Physical knowledge relationship.