**Figure 1: Shortcut learning examples used in QUEST performance evaluation stage. These synthetic images are extracted from the Image-Text pairs in the MS-COCO dataset. For the images, easily distinguishable MNIST number shortcuts are added, and for the text, corresponding number sequences to the images are included**

**Table 1: The results on the GTZAN and FMA-SMALL.**

| Method | GTZAN | | | | | | FMA | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Image to Audio | | | Audio to Image | | | Image to Audio | | | Audio to Image | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| InfoNCE | 34.01 | 84.73 | 94.41 | 32.48 | 78.68 | 90.86 | 15.87 | 28.62 | 35.87 | 12.50 | 25.50 | 29.12 |
| QUEST | **41.62** | **88.83** | **97.65** | **35.53** | **82.23** | **93.40** | **17.83** | **30.87** | **38.50** | **13.50** | **26.52** | **30.62** |

**Table 2: The results on the CLOTHO and AUDIOCAPS.**

| Method | CLOTHO | | | | | | AUDIOCAPS | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Text to Audio | | | Audio to Text | | | Text to Audio | | | Text to Image | | |
| | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 | R@1 | R@5 | R@10 |
| InfoNCE | 20.16 | 51.30 | 66.56 | 20.06 | 52.66 | 68.23 | 4.59 | 17.79 | 26.22 | 5.45 | 15.78 | 22.48 |
| QUEST | **21.10** | **52.45** | **68.86** | **22.36** | **54.23** | **70.42** | **5.16** | **18.08** | **27.17** | **6.02** | **15.98** | **24.78** |