



OUTLINE

- Introduction
- Related work
- Dataset
- Baseline
- Main approach
- Evaluation metric
- Results & analysis

Introduction

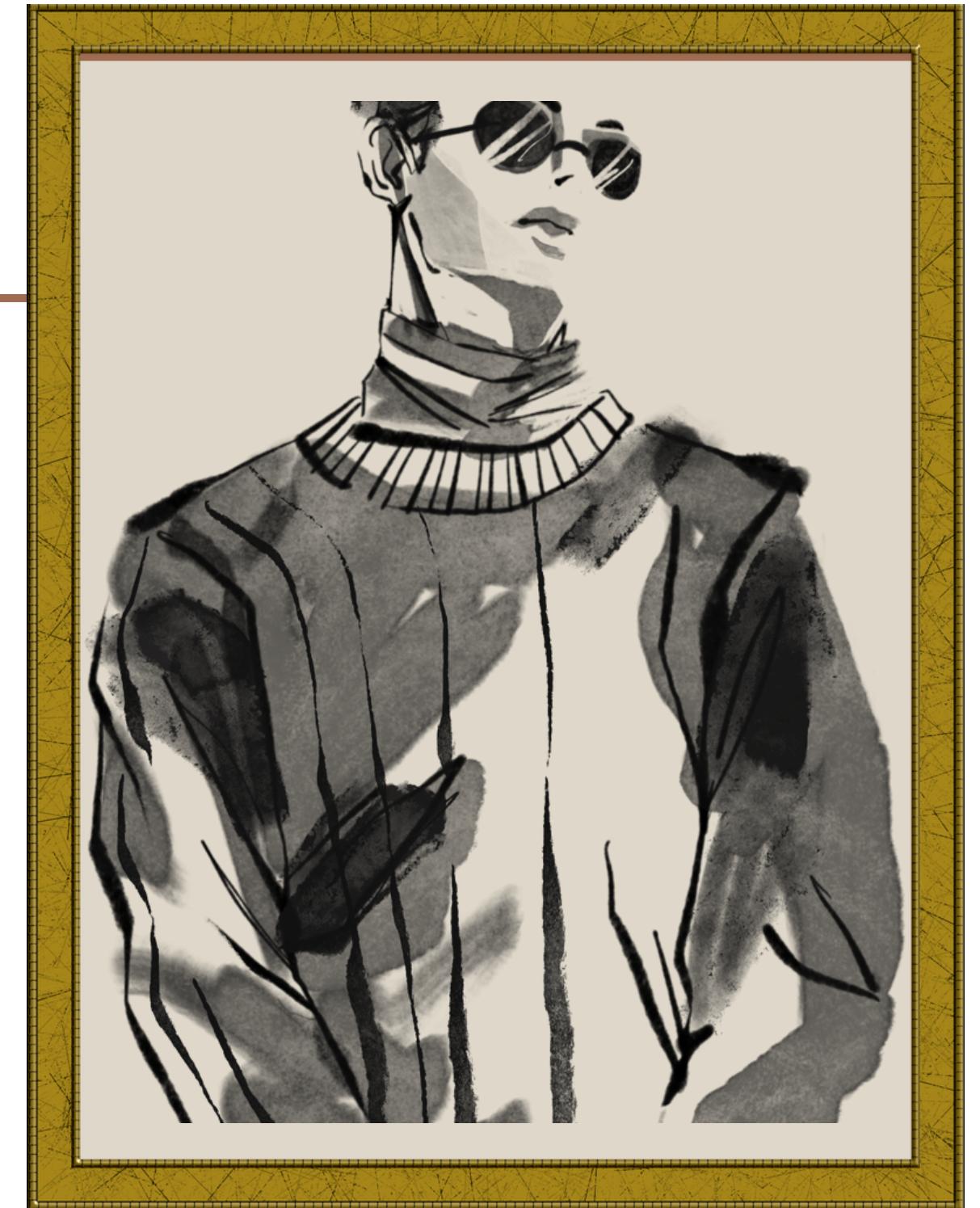
Why text to image?

Text to _____ ?

Intuitive? Interesting?

After taking the lesson, the question "text to ?" got imprinted in our mind. Therefore, when we are told to choose a topic for our final project, we decided to focus on this, and find the field we want to work on through two features: **Intuitive** and **Interesting**.

Intuitive means that the users can understand what our project is about easily. Interesting is that we wish the application can attract people to try on it. After discussing and brainstorming, we decided to work on text to image.



Introduction

Why dressing?

What do we care ?

What people of all ages care?

After deciding the field, we want our topic to be useful and attractive to all ages.

Therefore, we started thinking about what we ourselves need daily.

After listing things like how to set environment, how to make a presentation, how to get dressed, and a lot more, we find dressing the most common need for people.



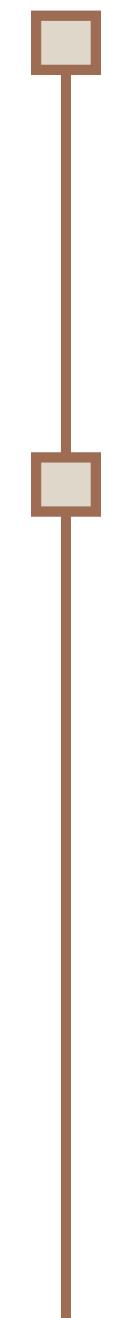
Related Work

HUMANGAN

- a VAE-based framework
- images generated by sampling from the learned distribution

OURS

- texture-aware codebook
- sampler with mixture-of-experts
- feedforward index prediction network for the hierarchical sampling
- allows for controllable human generation by giving texts describing the desired attributes



Datasets

- Downloaded from DeepFashion-MultiModal and extra datasets which we downloaded from another website and we annotated it **MANUALLY**
- Each image includes densepose and human parsing and attributes for both clothes shapes and fabric
- Contain 10425 train images and 1159 test images
- Size of the dataset: 750x1101



Datasets

Original Datasets

Contain normal clothing like T-shirt, and pants.

Extra Dataset

We made it more diverse (o^~^o) !!

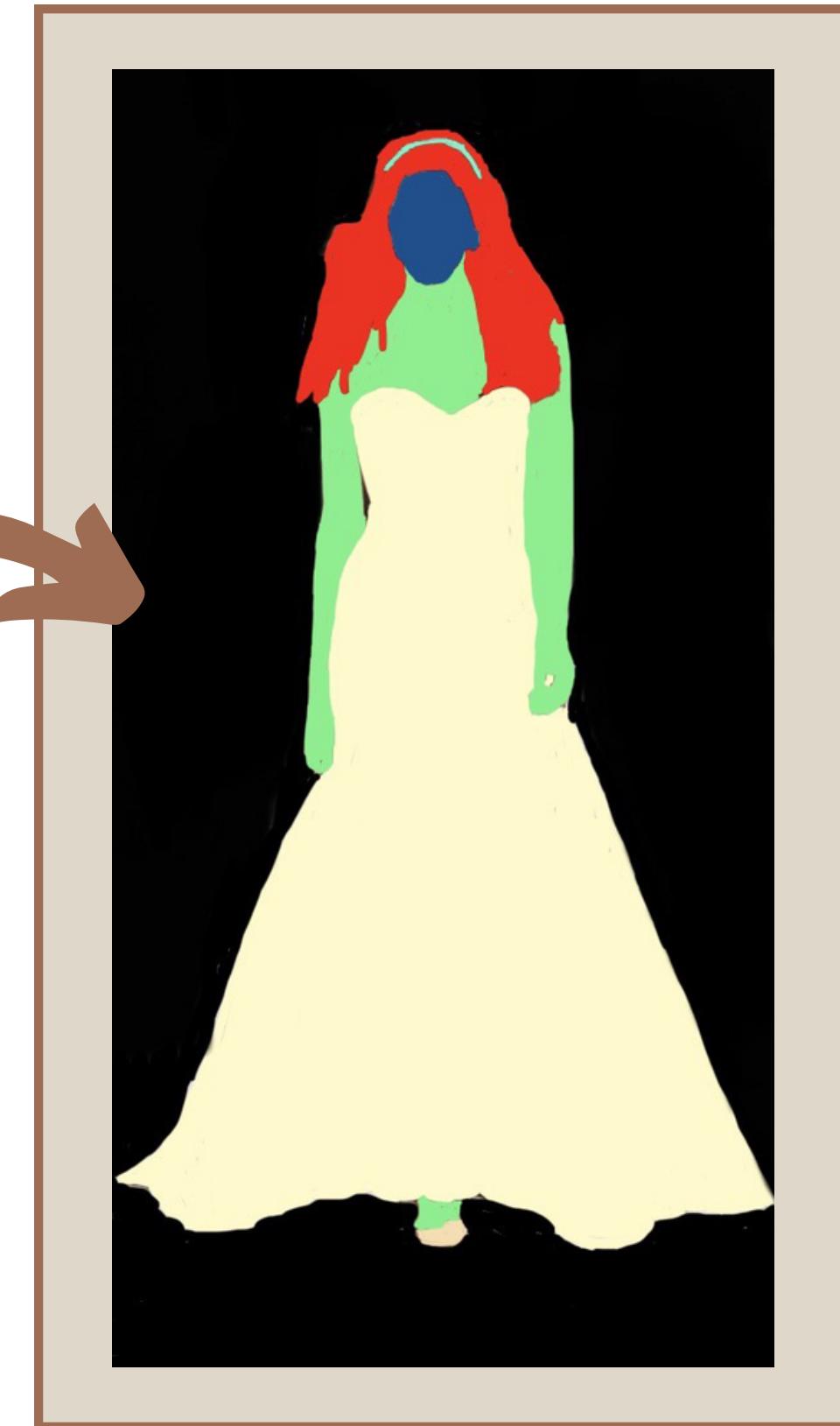


Datasets

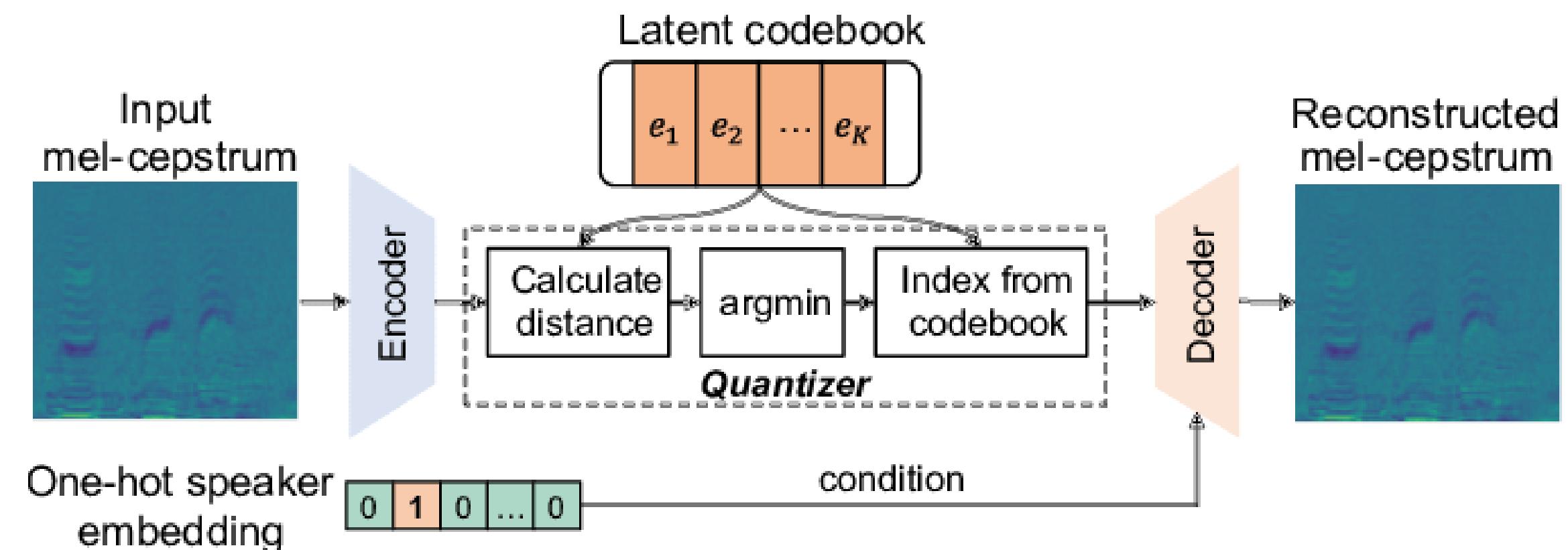
For the segmentation dataset, we drew the colors by hand. The colors could be recognized and each of them matches to one annotation

```
0: 'background' (0, 0, 0)
1: 'top' (255, 250, 250)
2: 'outer' (220, 220, 220)
3: 'skirt' (250, 235, 215)
4: 'dress' (255, 250, 205)
5: 'pants' (211, 211, 211)
6: 'leggings' (70, 130, 180)

7: 'headwear' (127, 255, 212)
8: 'eyeglass' (0, 100, 0)
9: 'neckwear' (50, 205, 50)
10: 'belt' (255, 255, 0)
11: 'footwear' (245, 222, 179)
12: 'bag' (255, 140, 0)
```

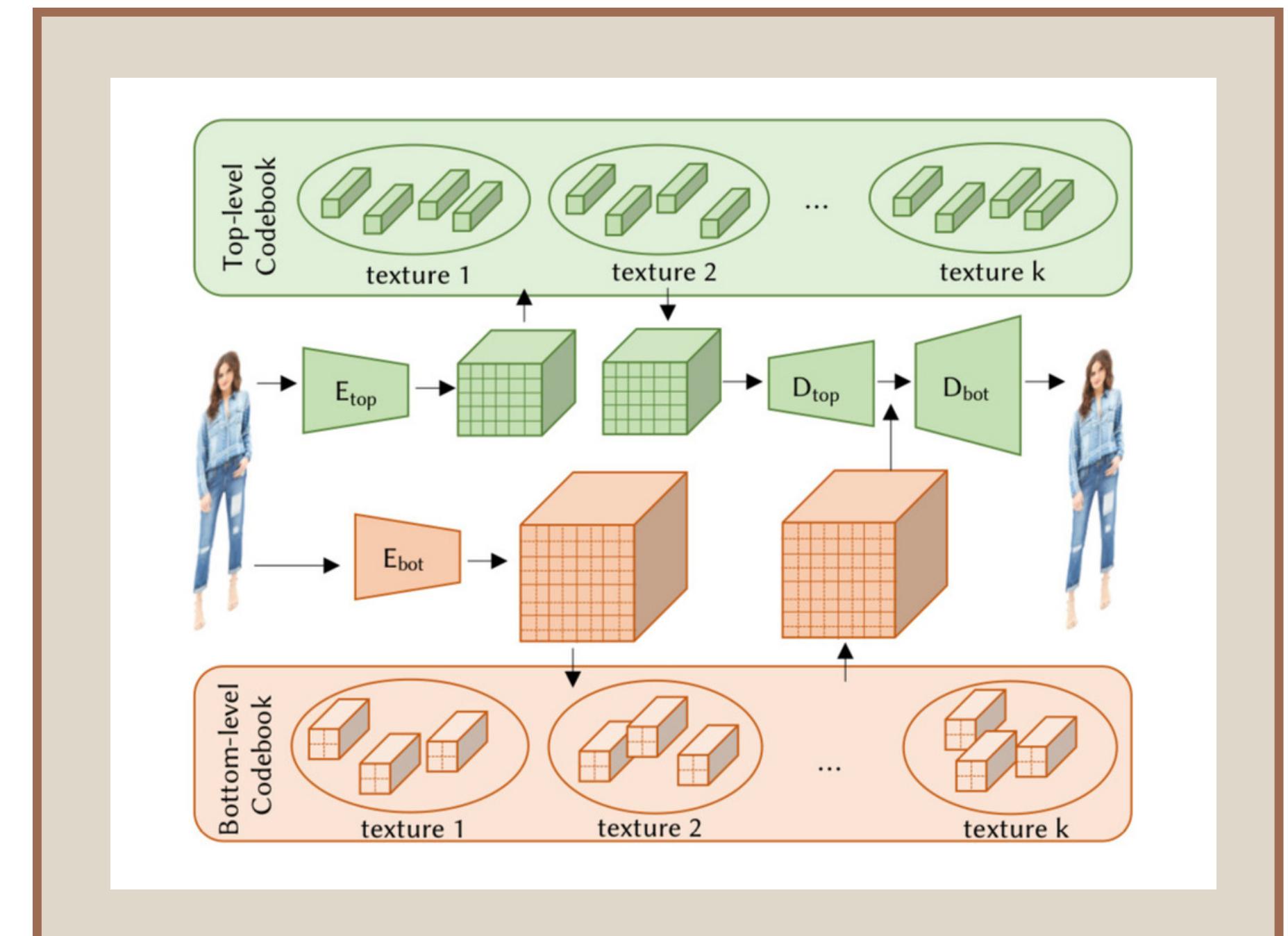


Baseline VQVAE



Baseline IMPLEMENTATION

- Is used to generate human images from the synthesized human parsing masks
- A hierarchical texture-aware codebook is built to store multi-scale neural representations for each type of texture
- The codebook indices are then generated using a diffusion-based transformer sampler with a mixture of experts
- To speed up the sampling process and refine the sampled images from the coarse level, a feed-forward codebook index prediction network is employed



Baseline

```
818 class Encoder(nn.Module):
819
820     def __init__(self,
821                  ch,
822                  num_res_blocks,
823                  attn_resolutions,
824                  in_channels,
825                  resolution,
826                  z_channels,
827                  ch_mult=(1, 2, 4, 8),
828                  dropout=0.0,
829                  resamp_with_conv=True,
830                  double_z=True):
831         super().__init__()
832         self.ch = ch
833         self.temb_ch = 0
834         self.num_resolutions = len(ch_mult)
835         self.num_res_blocks = num_res_blocks
836         self.resolution = resolution
837         self.in_channels = in_channels
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922     class Decoder(nn.Module):
923
924         def __init__(self,
925                     in_channels,
926                     resolution,
927                     z_channels,
928                     ch,
929                     out_ch,
930                     num_res_blocks,
931                     attn_resolutions,
932                     ch_mult=(1, 2, 4, 8),
933                     dropout=0.0,
934                     resamp_with_conv=True,
935                     give_pre_end=False):
936             super().__init__()
937             self.ch = ch
938             self.temb_ch = 0
939             self.num_resolutions = len(ch_mult)
940             self.num_res_blocks = num_res_blocks
941             self.resolution = resolution
942             self.in_channels = in_channels
943             self.give_pre_end = give_pre_end
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155     class Discriminator(nn.Module):
1156
1157         def __init__(self, nc, ndf, n_layers=3):
1158             super().__init__()
1159
1160             layers = [
1161                 nn.Conv2d(nc, ndf, kernel_size=4, stride=2, padding=1),
1162                 nn.LeakyReLU(0.2, True)
1163             ]
1164             ndf_mult = 1
1165             ndf_mult_prev = 1
1166             for n in range(1, n_layers): # gradually increase the number of filters
1167                 ndf_mult_prev = ndf_mult
1168                 ndf_mult = min(2**n, 8)
1169                 layers += [
1170                     nn.Conv2d(
1171                         ndf * ndf_mult_prev,
1172                         ndf * ndf_mult,
1173                         kernel_size=4,
1174                         stride=2,
1175                         padding=1,
1176                         bias=False),
1177                         nn.BatchNorm2d(ndf * ndf_mult),
1178                         nn.LeakyReLU(0.2, True)
1179                 ]
1180             ]
```

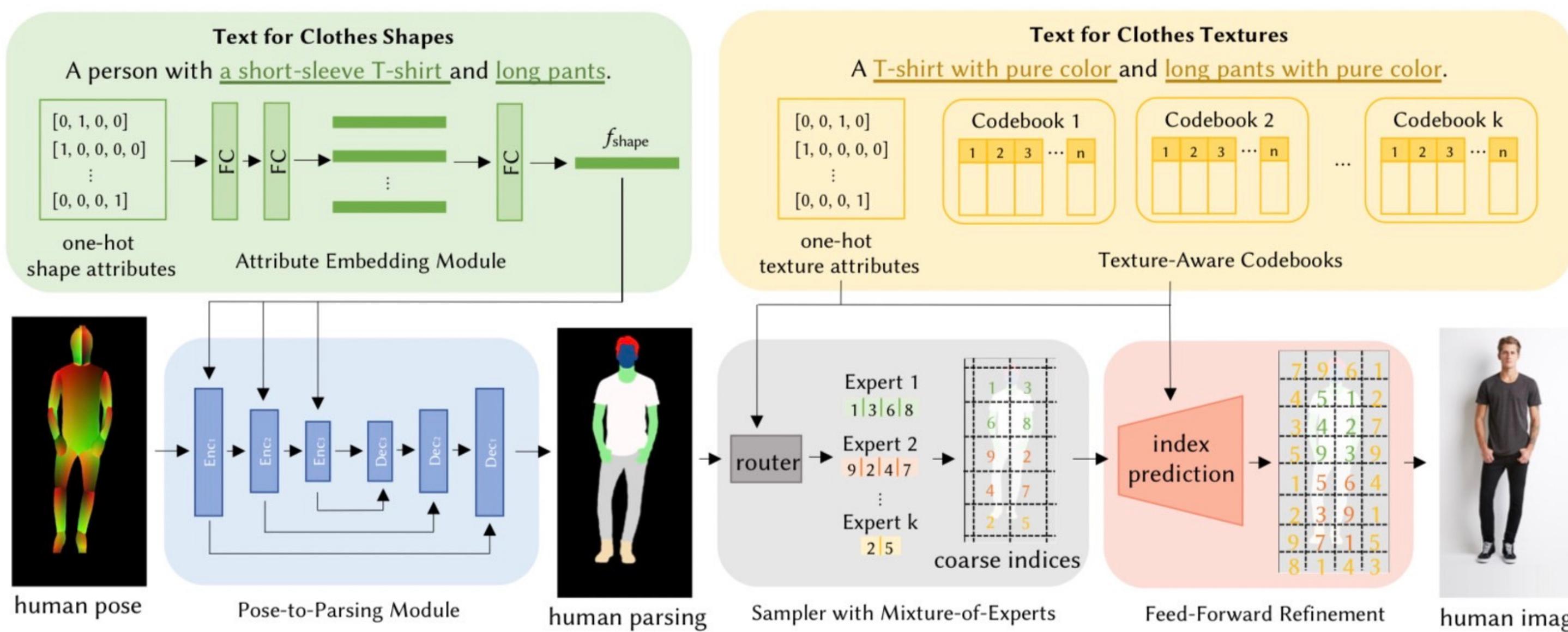
Baseline

Main Approach

We decomposed it into Two stages

- Stage I: Pose to Parsing
- Stage II: Parsing to Human

- Input: Densepose, and texts of shape and texture
- Output: A person wearing clothes (>^w^<)



Main Approach

STAGE 1: POSE TO PARSING

In the first stage of the project, a module called the Pose-to-Parsing Module was developed to generate a human parsing mask. This was done by combining human pose features with text descriptions of clothes' shapes. The text information was transformed into clothes **shape attributes**, which were then embedded and concatenated with the pose features. The encoder in the module spatially modulated the pose features and embedded attributes to produce the desired human parsing mask.

Main Approach

```
class SampleFromPoseModel(BaseSampleModel):
    """SampleFromPose model.

    """

    def __init__(self, opt):
        super().__init__(opt)
        # pose-to-parsing
        self.shape_attr_embedder = ShapeAttrEmbedding(
            dim=opt['shape_embedder_dim'],
            out_dim=opt['shape_embedder_out_dim'],
            cls_num_list=opt['shape_attr_class_num']).to(self.device)
        self.shape_parsing_encoder = ShapeUNet(
            in_channels=opt['shape_encoder_in_channels']).to(self.device)
        self.shape_parsing_decoder = FCNHead(
            in_channels=opt['shape_fc_in_channels'],
            in_index=opt['shape_fc_in_index'],
            channels=opt['shape_fc_channels'],
            num_convs=opt['shape_fc_num_convs'],
            concat_input=opt['shape_fc_concat_input'],
            dropout_ratio=opt['shape_fc_dropout_ratio'],
            num_classes=opt['shape_fc_num_classes'],
            align_corners=opt['shape_fc_align_corners'],
            ).to(self.device)
        self.load_shape_generation_models()

    def generate_parsing_map(self):
        with torch.no_grad():
            attr_embedding = self.shape_attr_embedder(self.shape_attr)
            pose_enc = self.shape_parsing_encoder(self.pose, attr_embedding)
            seg_logits = self.shape_parsing_decoder(pose_enc)
            self.segm = seg_logits.argmax(dim=1)
            self.segm = self.segm.unsqueeze(1)
```

Main Approach

STAGE 2: PARSING TO HUMAN

In Stage II, the goal is to generate a human image using the human parsing mask obtained in Stage I. This involves enhancing the mask with diverse textures of clothes based on user-provided texts. A hierarchical texture-aware codebook is used to learn and store neural representations of different textures. The human image generation process includes sampling indices from the codebook at multiple levels. Coarse-level indices are predicted using a mixture-of-experts sampler, while fine-level indices are refined using a feed-forward network. The aim is to create a realistic human image by incorporating a range of textures into the parsing mask.

Main Approach

```
def sample_and_refine(self, save_dir=None, img_name=None):
    # sample 32x16 features indices
    sampled_top_indices_list = self.sample_fn(
        temp=1, sample_steps=self.sample_steps)

    for sample_idx in range(self.batch_size):
        sample_indices = [
            sampled_indices_cur[sample_idx:sample_idx + 1]
            for sampled_indices_cur in sampled_top_indices_list
        ]
        top_quant = self.top_quantize.get_codebook_entry(
            sample_indices, self.texture_mask[sample_idx:sample_idx + 1],
            (sample_indices[0].size(0), self.shape[0], self.shape[1],
             self.opt["top_z_channels"]))
        top_quant = self.top_post_quant_conv(top_quant)

        bot_indices_list = self.bot_index_prediction(
            top_quant, self.texture_mask[sample_idx:sample_idx + 1])

        quant_bot = self.bot_quantize.get_codebook_entry(
            bot_indices_list, self.texture_mask[sample_idx:sample_idx + 1],
            (bot_indices_list[0].size(0), bot_indices_list[0].size(1),
             bot_indices_list[0].size(2),
             self.opt["bot_z_channels"]))#.permute(0, 3, 1, 2)
        quant_bot = self.bot_post_quant_conv(quant_bot)
        bot_dec_res = self.bot_decoder_res(quant_bot)
```

```
dec = self.decoder(top_quant, bot_h=bot_dec_res)
dec = ((dec + 1) / 2)
dec = dec.clamp_(0, 1)
if save_dir is None and img_name is None:
    return dec
else:
    save_image(
        dec,
        f'{save_dir}/{img_name[sample_idx]}',
        nrow=1,
        padding=4)
```

```
def generate_texture_map(self):
    upper_cls = [1., 4.]
    lower_cls = [3., 5., 21.]
    outer_cls = [2.]

    mask_batch = []
    for idx in range(self.batch_size):
        mask = torch.zeros_like(self.segm[idx])
        upper_fused_attr = self.upper_fused_attr[idx]
        lower_fused_attr = self.lower_fused_attr[idx]
        outer_fused_attr = self.outer_fused_attr[idx]
        if upper_fused_attr != 17:
            for cls in upper_cls:
                mask[self.segm[idx] == cls] = upper_fused_attr + 1
        if lower_fused_attr != 17:
            for cls in lower_cls:
                mask[self.segm[idx] == cls] = lower_fused_attr + 1
        if outer_fused_attr != 17:
            for cls in outer_cls:
                mask[self.segm[idx] == cls] = outer_fused_attr + 1
        mask_batch.append(mask)
    self.texture_mask = torch.stack(mask_batch, dim=0).to(torch.float32)
```

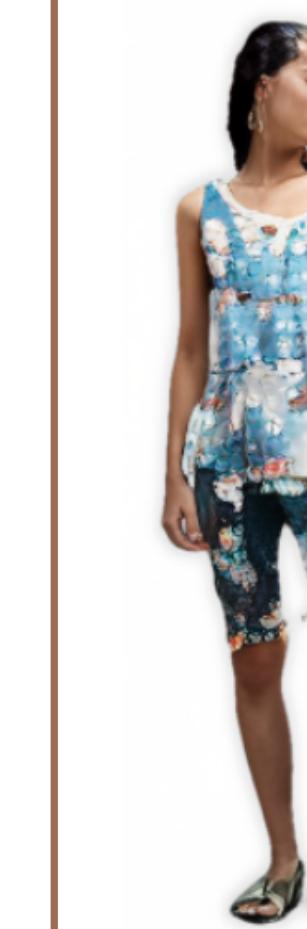
EXPERIMENT

Result and Analysis

- Experiment I: Different Sample Steps
- Experiment II: Different trained models, one with original dataset, another with extended dataset

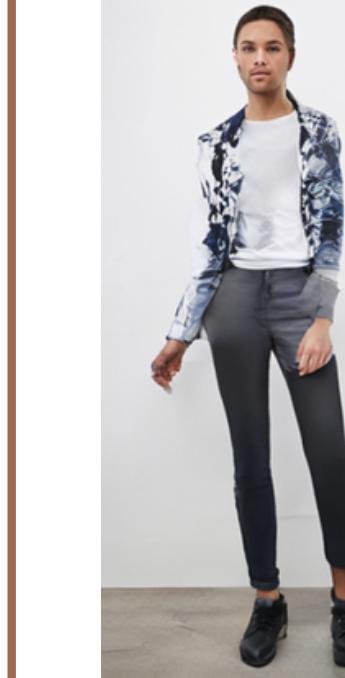


Evaluation Matric DIFFERENT SAMPLE STEPS

Sample Steps	50	100	150	200	250	300
Result						

We found that when sample steps is around to 250, it has the best result.

Evaluation Matric DIFFERENT MODELS

Original trained model					
Extended trained model					

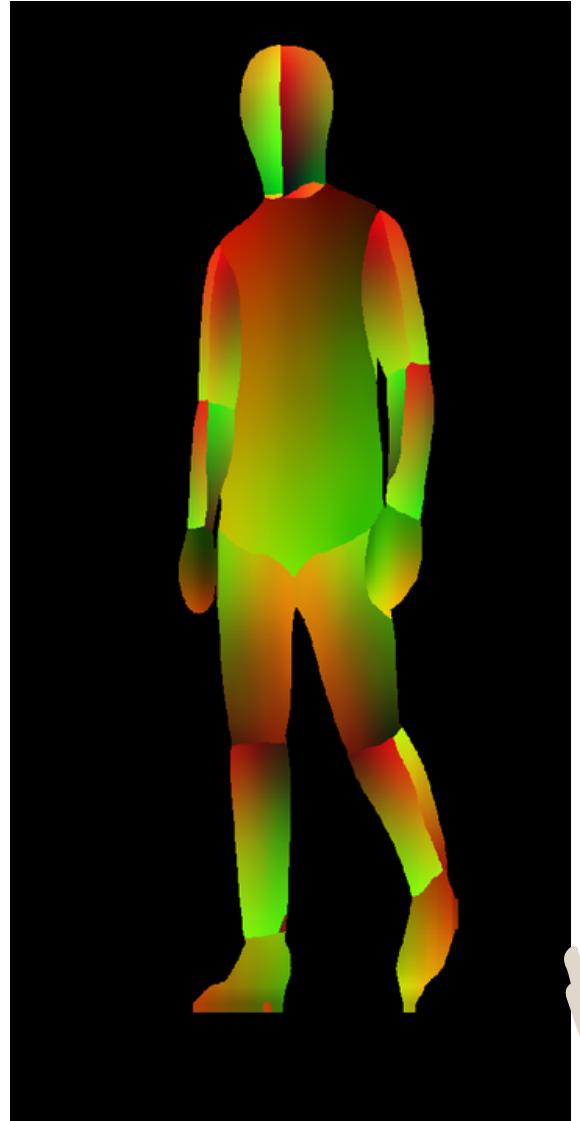
The below one is more diversity

Evaluation Matric DIFFFERENT TYPES OF CLOTHING

Short-pants	Long-pants	short-skirts	long-skirts	Short-jacket	Long-jacket
 (5,3.25)	 (4.75,3.75)	 (4,2.25)	 (1,1.25)	 (2.5,2.5)	 (3.75,4)
 (5,3.5)	 (5,4.25)	 (4.75,4.25)	 (3,3)	 (1.25,1.5)	 (3.25,3.25)
 (3.75,3)	 (4.25,3.5)	 (3.75,3)	 (2.5,2.25)	 (2,2.25)	 (2.75,2.5)

We give rates to each images by completeness and preference ,range from 1~5. The numbers below the images are the average rates given by our 4 members.

Pants perform better than skirts and jackets in completeness.



DISCUSSION AND ANALYSIS

Result and Analysis

We found that not every denseposes can generate the clothing with good quantity

Example: This densepose doesn't have a good quantity of generating long-pants



LIMITATION

Result and Analysis

We found that if the legs of the pose are not joined, the generated image would have broken skirt. The reason might be the bias caused by the dataset.



APPLICATION

Result and Analysis



We applied the result images to 3D images, giving us more information about dressing