

| Method | Pros | Cons |
|---|--|---|
| Pretrained CLIP (on b-boxes from object detection) | <ul style="list-style-type: none"> -No data required -No training required | <ul style="list-style-type: none"> -Most design algorithm for ranking based on b-boxes -Probably less accurate than trained models |
| Self-Supervised (e.g., SimCLR, EsViT) | <ul style="list-style-type: none"> -No labeled data required -Probably more accurate than untrained models | <ul style="list-style-type: none"> -Training required -Might capture features that are not related to clothes, because of the lack of label (can be used on b-boxes to fix the issue) |
| Classification Network | | <ul style="list-style-type: none"> -Training required -Labeled data required |
| Adding Extra Layer to Pretrained CLIP | <ul style="list-style-type: none"> -Probably great performance on the features that are included in captions | <ul style="list-style-type: none"> -Caption required for images (can be extracted from object detection features) -Only considers features that are included in captions for measuring similarity |