

# Controllable Recommenders using Deep Generative Models and Disentanglement

Samarth Bhargav

s.bhargav@uva.nl

IRLab, University of Amsterdam  
The Netherlands

Evangelos Kanoulas

e.kanoulas@uva.nl

IRLab, University of Amsterdam  
The Netherlands

## ABSTRACT

In this paper, we consider **controllability as a means to satisfy dynamic preferences of users**, enabling them to control recommendations such that their current preference is met. While deep models have shown improved performance for collaborative filtering, they are generally not amenable to fine grained control by a user, leading to the development of methods like deep language critiquing. We propose an alternate view, where instead of keyphrase based critiques, **a user is provided ‘knobs’ in a disentangled latent space, with each knob corresponding to an item aspect**. Disentanglement here refers to a latent space where generative factors (here, a preference towards an item category like genre) are captured *independently* in their respective dimensions, thereby enabling *predictable* manipulations, otherwise not possible in an entangled space. **We propose using a (semi-)supervised disentanglement objective** for this purpose, as well as **multiple metrics to evaluate the controllability and the degree of personalization of controlled recommendations**. We show that by updating the disentangled latent space based on user feedback, and by exploiting the generative nature of the recommender, controlled and *personalized* recommendations can be produced. Through experiments on two widely used collaborative filtering datasets, we demonstrate that a controllable recommender can be trained with a slight reduction in recommender performance, provided enough supervision is provided. The recommendations produced by these models appear to both conform to a user’s current preference and remain personalized.

## KEYWORDS

Collaborative Filtering, Representation Learning, Recommendation, Disentangled Representation Learning, Controllable Recommendation, Deep Critiquing

## 1 INTRODUCTION

Auto-encoder based architectures have shown impressive performance in collaborative filtering with implicit feedback [19, 25, 30]. However, these methods may fail to model short term or dynamic preferences of a user. Subsequently, methods to tackle this explicitly problem have been proposed, for example, through conversations via a conversational recommender system [13], or critiquing recommendations using language or keyphrases [18, 23, 35, 37]. In contrast to keyphrases, typically mined from reviews or descriptions (for instance), this work considers building controllable or critique-able recommenders using attribute data for items, which can be used to construct preference distributions at a user level. This is in turn used in the critiquing process or to express a short-term preference, allowing a user to exert control over a recommender in a meaningful manner. For instance, a user who typically watches

*Action* movies, or (as a consequence) exclusively receives *Action* recommendations, can now explicitly request *personalized* movies of other genres, like *Animation*.

Controllability can be achieved by utilizing the *generative* nature of certain *Deep Generative* Recommenders [18, 23, 35, 37]. Such models have an encoder which produces a user-latent representation, which is fed to the decoder to predict items for that user. A manipulation of this representation, followed by decoding step using this, should alter recommendations produced by the model. However, since the latent spaces of these models are typically *entangled*, the outputs produced through such manipulations are likely to be unpredictable or random, making them unusable for this task.

We propose using a *disentangled* latent space which, in contrast, can be manipulated predictably. A representation is disentangled w.r.t known ground truth variables or generative factors (ex. genres, availability, context, etc), if there is *only one* latent dimension in the representation that is influenced when this ground truth variable is changed [3, 12, 20, 21].

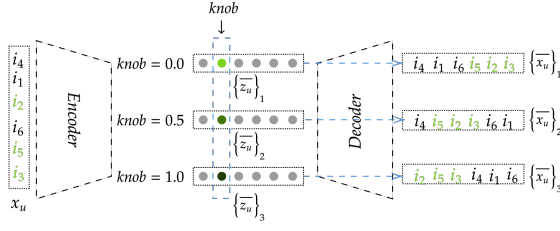
Prior work tackles critiquing/controllability by modeling a latent space where keyphrases are *co-embedded* with user/item embeddings. By ‘zero-ing’ out a certain keyphrase, the corresponding keyphrase embedding is used to update the user representation, producing critiqued recommendations. Our work in contrast doesn’t utilize a co-embedding space, and instead the latent space is directly manipulated. This means no additional optimization (like in [18, 22]) is required to incorporate (multi-step) critiques. Furthermore, this also allows for ‘positive’ critiques as well as ‘soft’ control (gradual, non-binary) instead of only binary critiques.

We propose using supervision to obtain a mapping of a particular factor/aspect to a latent dimension. Since supervision signals might not be available for all users, we experiment with settings where limited data is provided.

In addition, while there have been metrics previously proposed for evaluating controllability/critiquing [32, 35, 37], these metrics don’t explicitly account for the personalization of post-critique recommendations. We propose multiple metrics to evaluate the personalization of critiqued recommendations, measuring both binary critiques as well as ‘soft’, gradual controllability. To summarize, our contributions are as follows <sup>1</sup>:

- We propose using *(semi-)supervised disentanglement* to learn disentangled representations for controllable, *personalized* recommendations. Through experiments on two large scale collaborative filtering datasets using 2 types of signals as ‘knobs’, we show the proposed models produce controllable recommendations, at the cost of a slight reduction in overall recommendation performance. In addition, we experiment

<sup>1</sup><https://github.com/samarthbhargav/disentangled-control-recsys>



**Figure 1: Deep Controllable Recommendations with (semi-)supervised disentanglement:** A user’s interactions are captured in  $x_u$ , used to infer the user representation  $z_u$ . Given a latent dimension  $j$  (the ‘knob’), corresponding to a known factor (green items), only the values corresponding to the  $j$ th dimension of  $z_u$  are replaced to produce  $\{\tilde{z}_u\}_{1,2,3}$ . Each of these are then decoded to produce recommendations. The manipulation here pushes the green items higher in the ranked list, leading to controllable recommendations.

with different levels of supervision, and show controllability can be achieved with limited data.

- We focus on retaining recommendation performance while achieving effective controllability, and propose several metrics to extensively study (a) the degree of personalization of controlled recommendations (b) whether control is achieved in isolation (i.e if only one factor changes at a time) and (c) the effect of disentanglement on controllability. Using these metrics, we demonstrate that such models are amenable to user control, and the controlled recommendations appear to be personalized to an extent.

## 2 BACKGROUND AND RELATED WORK

### 2.1 Deep Recommender Systems

There have been several works that use Deep Learning for Recommendation [11, 31, 33, 38]. A common theme in several deep models is the use of auto-encoder architectures like the Collaborative Denoising Autoencoder (CDAE) [36] or MultVAE [19]. The latter uses the Variational Autoencoder framework [5, 17] for recommendation. Shenbin et al. [30] propose the state-of-the-art RecVAE. Deep Recommender systems however are (typically) black-box models which are difficult to interpret, compared to content-based methods [38].

### 2.2 Disentangled Representation Learning

The use of deep *generative* recommenders allow for disentangled representations, regaining some interpretability [9, 24, 34]. Most such models use the VAE framework, where the objective is to *reconstruct* the input with high probability, with an additional loss term *constraining the latent space*. By imposing additional constraints, *disentangled* representations can be learnt. These methods are typically unsupervised and therefore come with limitations [20]: they were found to be very sensitive to hyperparameters, *reliably* learning unsupervised disentangled representations is a very challenging task.

There have been several works that use disentanglement in recommendation: Ma et al. [25] assume that disentanglement is generated by user behaviours on ‘macro’ and ‘micro’ levels, and show that their models outperform non-disentangled baselines. Ma et al. [26] apply disentanglement to the sequential recommendation task, while Wang et al. [34] disentangle diverse user-intents using graph based collaborative filtering; Cui et al. [9] propose DGCF, which have ‘implicit’ (unknown) and ‘explicit’ (known) signals that are disentangled using an RNN based model with a two-step method since some computations are non-differentiable. Wang et al. [32] propose using weakly-supervised disentanglement objective on pairs of items. This allows for attribute-based item retrieval, such that items differ more/less on the provided attribute. In contrast to prior work, we use semi-supervised disentanglement, while focusing on using disentanglement for controllability.

### 2.3 Critiquing / Controllable Recommenders

We focus on critiquing for *deep* recommenders (for a full overview of earlier work in critiquing, see Chapter 13 of [29]). The following paragraphs detail approaches that leverage language or keyphrases for critiquing, followed by other approaches like ours which does not utilize any language/keyphrase. We end this sub-section with a discussion about evaluation metrics.

Recent work has focused on using language or keyphrase based critiquing [18, 23, 35, 37]. Wu et al. [35] propose the Deep Language based Critiquing (DLC) paradigm, where (subjective) keyphrases are treated both as explanations as well as a means for a user to critique recommendations. Here, critiquing is achieved by rejecting or ‘zero’-ing out keyphrases which are co-embedded with users. Luo et al. [22] adapt the CE-NCF for *multi-step* critiquing; [18] use a ranking optimization instead of pairwise re-scoring to re-rank items; [23] use a VAE-based model to achieve better training stability and lower computational complexity; [37] propose using Keyphrase Activation Vectors instead of using a second ‘head’ used in [35], and adapt it for ‘positive’ critiques. Antognini et al. [2] propose T-RECS, which infers keyphrases from the intersection of user profiles and an item. [1] show that a model under a weak-supervision scheme matches or beats recommendation/explanation and critiquing performance while being much faster. We use supervised disentanglement in the user-latent space instead of a separate embedding space, while not requiring any language data (at the cost of no explanations). Prior work focuses on binary signals ex. reject/accept, while our work considers ‘soft’ critiques, allowing for a *gradual* tuning. Note that while we experiment only with unit critiques, compound and/or multi-step critiques is also possible within our proposed models, which we leave for future work.

Wang et al. [32] propose an orthogonal task, where items are retrieved on a ‘gradient’ given an attribute, focusing on *gradual* change, similar to our work. Our work focuses on recommendation performance whereas [32] focus on gradient item retrieval. Cen et al. [7] propose ComiRec for sequential recommenders, where control is used to balance diversity and recommendation accuracy, from the perspective of a designer. The method presented in [25] is closest to us, where representations are first altered, followed by a nearest neighbours search over items using a beam search. In contrast, we utilize the *generative nature* of the model instead of

performing nearest neighbour search, in addition to using semi-supervised disentanglement. The authors are aware of concurrent work [28] that is close to ours. Since the paper/code was unavailable at the time of experiments/ publication, we leave the evaluation of [28] for future work.

*Metrics for critiquing/controllability.* Wu et al. [35] propose the Falling MAP for *negative* critiques. Given a user and a keyphrase (here, genre/tag) this metric measures the difference in MAP values computed on the set of items that have the critiqued keyphrase. If this value is positive, items of the critiqued keyphrase ‘falls’ down the ranking list, which is desirable. Similarly [37] propose a follow up metric, which accounts for positive critiques as well, comparing the normalized difference of average ratings of items before / after the critique. We note that these metrics don’t explicitly model the *personalization* of post-critique recommendations. In contrast, the proposed metrics explicitly consider this, by measuring performance on against items for which a user’s preference is known.

### 3 METHODOLOGY

We present disentanglement of the latent space using semi-supervision as a means to control in Section 3.1. Adapting disentanglement representation learning for this task is explored in detail in Section 3.2.

#### 3.1 Control using (semi-)supervised disentanglement

We train a model using a (semi-)supervised disentanglement loss, producing disentangled user representations. During inference, representations for all users are inferred (possibly offline). In the foreground, the representations are manipulated on the fly by increasing/decreasing a ‘knob’, allowing for granular change. In the background, the mapping of a knob to a latent dimension is utilized to modify the user’s representation, which is fed to the decoder to produce controlled recommendations. The following paragraphs motivate the use of disentanglement and semi-supervision to build such a recommender.

As mentioned before, the entangled nature of a typical deep recommender renders meaningful manipulation impossible. This is because (a) there is no guarantee that manipulations produce changes in a single generative factor, and (b) it is not apparent *which* dimension to alter given a factor. Even if we assume (b), predictability is not guaranteed, ex. manipulating the *Thriller* dimension may increase the number of *Thriller* movies, but it may inadvertently increase the number of *Children* movies, which is undesirable.

We propose using supervised disentanglement to tackle these issues. Disentanglement can isolate a generative factor to a single dimension, and altering this dimension (assuming perfect disentanglement) space should produce a change only in this generative factor. Using (semi-)supervised disentanglement addresses two issues: (a) unsupervised disentanglement has several shortcomings [20], which semi-supervision may alleviate; and (b) a particular dimension is now ‘tied’ to a generative factor, ensuring that the correct dimension is being altered.

These recommendations should ideally be *personalized*, since they were trained with a reconstruction loss. In essence, a recommender system can be tuned *on the fly* to express user preferences for the current session. This process is illustrated in Figure 1.

In this paper we opt for preference distributions i.e tag/genre distributions, as generative factors or ‘knobs’ (see Section 5.1.2), since they are widely available. However, preference distributions are *not independent* and can also be very *noisy* [4, 10], in contrast to independent and exact generative factors used in disentanglement literature e.g dSprites[27]. Consequently, they may be difficult to learn and/or disentangle. The following section discusses particulars about the models.

#### 3.2 Unsupervised and Semi-supervised Disentanglement

In this section, two models, the  $\beta$ -VAE [12, 19] and the  $\beta$ -TCVAE [8], are adapted for the recommendation task using a Multinomial loss [19], resulting in recommenders that can produce disentangled representations. Furthermore, this objective is further supplemented with a semi-supervision loss.

**3.2.1 Unsupervised Disentanglement.** The VAE [17] optimizes the marginal (log-)likelihood of the observed data  $\mathbf{x}$  in expectation over the whole distribution of latent factors  $\mathbf{z}$ . It assumes that the data is supported on a low dimensional manifold in a high dimensional space. The assumption of a latent code allows us to express the marginal distribution as follows:  $p_{\theta}(\mathbf{x}) = \int_{\mathbf{z}} p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ . Since this quantity is intractable, the Evidence Lower Bound is optimized instead [14]:

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}_{VAE} = \mathbb{E}_{\mathbf{x}}[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - KL[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]] \quad (1)$$

where  $q_{\phi}(\mathbf{z}|\mathbf{x})$  is a variational approximation of the posterior distribution, which is typically a Gaussian distribution.  $p(\mathbf{z})$ , the prior, is usually an isotropic Gaussian:  $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Several methods modify this objective to encourage disentanglement:

**3.2.2  $\beta$ -VAE.** The  $\beta$ -VAE[6, 12] modifies the VAE loss i.e Equation 1 by adding a  $\beta$ -multiplier to the KL term:

$$\mathcal{L}_{\beta VAE} = \mathbb{E}_{\mathbf{x}}[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \beta.KL[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]] \quad (2)$$

Disentanglement is encouraged by setting  $\beta > 1$ , which further forces the posterior distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$  to be close to the *isotropic* prior for every  $\mathbf{x}$  (as opposed to only on average), imposing statistical independence. However, this involves a trade-off between reconstruction quality and disentanglement [6, 12]. While Liang et al. [19] use  $\beta < 1$  for better recommender performance, we opt for  $\beta \geq 1$  to encourage disentanglement, and achieve similar/better recommender performance by removing the Dropout layer, following Shenbin et al. [30]. An annealing strategy is used for  $\beta$ , following [6, 19]. The resulting model, which uses a Multinomial loss (see Section 5.3.2), is termed Mult- $\beta$ -VAE.

**3.2.3  $\beta$ -TCVAE.** Chen et al. [8] show that by imposing a constraint on only a part of the KL divergence term, reconstruction and disentanglement increases compared to the  $\beta$ -VAE model. The KL term

is decomposed into 3 terms (a) an index-code Mutual Information (MI) term, (b) dimension-wise KL, and (c) a total correlation (TC) term. The index-code MI captures the information between the data and latent variables, based on the empirical data distribution; the dimension-wise KL prevents deviation of the posterior from the prior dimensions; and the TC term encourages statistical independence of the posterior distribution leading to disentanglement. The key idea is that the  $\beta$ -VAE constrains all 3 terms, which might harm performance, while only the TC term is penalized in this model. The loss function of the TC-VAE is as follows:

$$\begin{aligned} L_{\beta\text{-TCVAE}} = & \mathbb{E}_{\mathbf{z}} [\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \alpha I_q(\mathbf{z}; \mathbf{n}) \\ & - \gamma \sum_j KL[q(\mathbf{z}_j) || p(\mathbf{z}_j)] - \beta KL[q_{\phi}(\mathbf{z}) || \prod_j q_{\phi}(\mathbf{z}_j)]] \end{aligned} \quad (3)$$

where  $\alpha$ ,  $\gamma$  and  $\beta$  are multipliers for the index code MI, the dimension wise KL and the TC term respectively. Following Chen et al. [8],  $\gamma = \alpha = 1$  is used,  $\beta$  is a hyperparameter. The final model, which we term Mult-TCVAE, uses a Multinomial loss (see Section 5.3.2).

**3.2.4 Semi-supervised Disentanglement.** The distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$  can be constrained so that a given dimension of  $q(\mathbf{z})$  is associated with a ground truth attribute. That is, given a representation (for instance, the mean vector from the posterior) for a data point,  $\boldsymbol{\mu} \in \mathbb{R}^D$ , and the set of known attributes  $\mathbf{a} \in \mathbb{R}^A$  for that data point, with  $A \leq D$ , the  $j$ th dimension of  $\boldsymbol{\mu}$  is predictive of the  $j$ th dimension of  $\mathbf{a}$ ,  $1 \leq j \leq A$ . This is achieved by adding a semi-supervised loss  $R_s$ :

$$\mathcal{L}_{ss} = \mathcal{L}_{\text{unsup}} + \gamma_{ss} \cdot \mathbb{E}_{\mathbf{x}, \mathbf{z}} R_s(q_{\phi}(\mathbf{z}|\mathbf{x}), \mathbf{a}) \quad (4)$$

$\mathcal{L}_{\text{unsup}}$  is either Equation 2 or 3. The binary cross entropy loss is used for  $R_s$  [21]:

$$R_s(\boldsymbol{\mu}, \mathbf{a}) = \sum_{i=1}^A \mathbf{a}_i \log(\sigma(\boldsymbol{\mu}_i)) + (1 - \mathbf{a}_i) \log(1 - \sigma(\boldsymbol{\mu}_i)) \quad (5)$$

where  $\sigma(\cdot)$  is the logistic function. Note that the remaining dimensions of  $\mathbf{z}$ , those that are not supervised, are not constrained by the semi-supervision loss.

## 4 EVALUATION METHODOLOGY

In this section we propose an evaluation methodology and a set of metrics that investigates the extent of controllability and personalization of a recommender, described in Section 4.1. We also briefly describe existing metrics for evaluating disentanglement in Section 4.2.

### 4.1 Evaluating Controllability

For ease of discussion, this section assumes a single factor  $g_j$ , s.t given a user  $u$ , and a factor  $g_j$ , the set of all items with factor  $g_j$  is  $\mathcal{I}_{g_j}$ , items encountered/rated by a user is  $\mathcal{I}_u$ , and items with factor  $g_j$  that user  $u$  has encountered is  $\mathcal{I}_{(u,g_j)} = \mathcal{I}_{g_j} \cap \mathcal{I}_u$  (subscripts for  $\mathbf{z}$  are dropped for brevity in the following descriptions). Note that

**Table 1: Metrics for Controllability, listed with the inputs and holdouts used in its computation.  $g_j$  is the genre being manipulated, while  $g_i$  is another genre**

Metric	Input	Holdout(s)
$\delta_{\text{ctrl}}$	$\mathcal{I}_u - \mathcal{I}_{(u,g_j)}$	$\mathcal{I}_{(u,g_j)}$
$\delta_{\text{irrel}}$	$\mathcal{I}_u - \mathcal{I}_{(u,g_j)}$	$\mathcal{I}_{g_j} - \mathcal{I}_{(u,g_j)}$
Corr	$\mathcal{I}_u - \mathcal{I}_{(u,g_j)}$	$\mathcal{I}_{(u,g_j)}$
Corr <sub>ctrl</sub>	$\mathcal{I}_u - \mathcal{I}_{(u,g_j)} - \mathcal{I}_{(u,g_i)}$	$\mathcal{I}_{(u,g_j)}$
Corr <sub>rand</sub>	$\mathcal{I}_u - \mathcal{I}_{(u,g_j)} - \mathcal{I}_{(u,g_i)}$	$\mathcal{I}_{(u,g_i)}$

an item can have multiple factors e.g a movie can be both *Action* and *Adventure*.

Given a user representation  $\mathbf{z}$ , the dimension corresponding to  $g_j$  is obtained by using a mapping function  $k$ , s.t  $k(g_j)$  is the dimension corresponding to  $g_j$  in  $\mathbf{z}$ , with  $k(g_j) \in \{1, \dots, K\}$ , with  $K \leq D$ , where  $D$  is the dimensionality of the latent space. In addition, we assume that manipulations of each dimension are in  $\in [0, 1]$  (0 is minimum). A manipulation entails the replacement of the  $k(g_j)$  dimension with a value indicating the position of a 'knob', producing  $\tilde{\mathbf{z}}$ . In addition, the top- $n$  items of  $\tilde{\mathbf{x}}$  is denoted by  $\text{top}_n(\tilde{\mathbf{x}})$ , and the number of items of a factor  $g_j$  in the top- $n$  list is denoted by  $\text{Count}(\text{top}_n(\tilde{\mathbf{x}}), g_j)$ . Finally, the  $i^{\text{th}}$  element of a vector  $\mathbf{a}$  is denoted by  $[\mathbf{a}]_i$ , with  $[\mathbf{b}]_i = 0$  indicating that the  $i^{\text{th}}$  dimension is set to 0.

**4.1.1 Desiderata.** For evaluating controllability, there are a list of desiderata:

- (1) Increasing the values of the  $k(g_j)$  dimension should push items of  $g_j$  higher. In other words, given  $[\tilde{\mathbf{z}}^m]_{k(g_j)} > [\tilde{\mathbf{z}}^n]_{k(g_j)}$  and assuming  $\tilde{\mathbf{z}}^m$  produces  $\tilde{\mathbf{x}}^m$  and  $\tilde{\mathbf{z}}^n$  produces  $\tilde{\mathbf{x}}^n$  respectively, then:

$$\text{Count}(\text{top}_n(\tilde{\mathbf{x}}^m), g_j) > \text{Count}(\text{top}_n(\tilde{\mathbf{x}}^n), g_j)$$

for small values of  $n$ . A natural consequence of this is irrelevant items  $\mathcal{I}_{g-j}$  should be pushed down or be replaced with relevant items, i.e.  $\text{Count}(\text{top}_n(\tilde{\mathbf{x}}^m), g-j)$  should decrease. Put another way,  $k(g_j)$  should only control items belonging to  $g_j$ , and *not* control items of other factors.

- (2) In addition, it is crucial to ensure that these recommendations are *personalized* i.e  $\text{top}_n(\tilde{\mathbf{x}}^m)$  should have items of the requested genre, but only those items that the user might like, i.e. items in  $\mathcal{I}_{(u,g_j)}$ .

We assume a ranking metric R-MET, computed against a held out set of items. R-MET can be any ranking metric, e.g. NDCG, assuming the relevance of all items except the ones in the holdout set to zero. We consider five metrics in total, with the two  $\delta$  metrics,  $\delta_{\text{ctrl}}$  and  $\delta_{\text{irrel}}$ , measuring the ranking changes of items of interest versus other items; and three correlation metrics measuring the gradual change observed as a dimension is manipulated: Corr, and (Corr<sub>ctrl</sub>, Corr<sub>rand</sub>). The  $\delta$  metrics measures the change induced when a user prioritizes items of only one factor, reflecting a critique or a short term preference i.e 'show me *Animation* movies only'. The correlation metrics on the other hand can capture exploratory behaviours or soft preferences ('maybe I will try at *Horror* movies'), and measures to what extent the recommender is *interactive* and *predictable*. These metrics are summarized in Table 1 along with

their inputs and holdout sets. Note that all the metrics are averaged across genres and then users.

**4.1.2  $\delta$  metrics.** Similar to *F-MAP* [35] and *PostCritRatingDiff* [37],  $\delta_{\text{ctrl}}$  measures the change in the ranking as a knob is set to its maximum value. Note, however, that only items *not* belonging to  $g_j$  that the user likes are fed into the model, with the intent to observe if manipulating the  $k(g_j)^{\text{th}}$  dimension increases the positions of the held-out items of genre  $g_j$  that a user likes. As a consequence, the model has to infer the correct items of  $g_j$  to recommend, without knowledge of the users preference towards items of  $g_j$ . Furthermore, in contrast to prior metrics, this evaluation is *personalized* i.e the metric is computed against  $I_{(u,g_j)}$  and not  $I_{g_j}$ . The process to compute the metric is detailed below.

First, *before* the manipulation, items of *other* genres that the user has liked,  $I_u - I_{(u,g_j)}$  is used to infer  $\mathbf{x}$ .  $\text{R-MET}_{\text{default}}$  is computed on this, using  $I_{(u,g_j)}$  as the holdout set, measuring the baseline ranking score of *personalized* items belonging to  $g_j$ . Second,  $[\bar{\mathbf{z}}]_{k(g_j)}$  is set to max indicating a high preference, and decoded to produce  $\bar{\mathbf{x}}$ . This is used to compute  $\text{R-MET}_{\text{ctrl}}$  on the same holdout, expressing the ranking score of *personalized* items belonging to  $g_j$  *after* manipulation. The difference between the two metrics is  $\delta_{\text{ctrl}}$ :

$$\delta_{\text{ctrl}} = \text{R-MET}_{\text{ctrl}} - \text{R-MET}_{\text{default}} \quad (6)$$

The higher the  $\delta_{\text{ctrl}}$ , the better the ability to produce personalized and controlled recommendations.

The second metric,  $\delta_{\text{irrel}}$ , is a variation of  $\delta_{\text{ctrl}}$ . While recommending items of genre  $g_j$  before items of other genres is important, it is desirable to recommend *items that the user might like*, over irrelevant items of  $g_j$ .  $\delta_{\text{irrel}}$  is similar to  $\delta_{\text{ctrl}}$ , except the holdout set are items of  $g_j$  *not* rated by the user,  $I_{g_j} - I_{(u,g_j)}$ . If  $\delta_{\text{irrel}}$  is high, items that the user *might not like* are also being promoted, which may be undesirable. While these two metrics reflect the extent of controllability, granular metrics, reflecting *gradual* change is presented in the following section.

**4.1.3 Correlation Metrics.** The following metrics compute the correlation between changes in ranking metric R-MET and the *values* of the manipulated dimension. The correlation measure used in the experiments is the Pearson’s correlation, we leave non-linear variants to future work.

**Corr** quantifies the correlation of *gradual* manipulations (as opposed to setting it to the max value as with  $\delta_{\text{ctrl}}$ ) with R-MET using  $I_{(u,g_j)}$  as the holdout set. **Corr<sub>ctrl</sub>** and **Corr<sub>rand</sub>** measure the same correlation with different holdout sets and identical inputs. **Corr<sub>ctrl</sub>** uses  $I_{(u,g_j)}$  as the hold out, measuring the effect of manipulation against items of  $g_j$ , whereas **Corr<sub>rand</sub>** uses  $I_{(u,g_i)}$  which measures the effect of manipulation against items of a random *control* factor  $g_i$ . The input set for the last 2 metrics has neither items of  $g_j$  nor  $g_i$ . In the ideal case,  $\text{Corr}_{\text{ctrl}}$  should be positive i.e controlling  $k(g_j)$  positively influences items of  $g_j$ , but  $\text{Corr}_{\text{rand}}$  should be zero (no change) or negative (decrease, either due to replacement or demotion), meaning items of another factor are not/negatively influenced.

It is our hypothesis that keeping items of irrelevant genres low in the ranking is harder if the genre being controlled is highly correlated with another genre in the data (e.g many *Action* movies

are also tagged *Adventure*). To measure this, we use two  $g_i$ :  $g_{\text{easy}}$  and  $g_{\text{difficult}}$ , where  $g_{\text{easy}}$  is the least co-occurring factor with  $g_j$  and  $g_{\text{difficult}}$  conversely frequently co-occurs with  $g_j$  producing  $\text{EasyCorr}_{\text{ctrl}}$ ,  $\text{EasyCorr}_{\text{rand}}$ ,  $\text{DiffCorr}_{\text{ctrl}}$  and  $\text{DiffCorr}_{\text{rand}}$ <sup>2</sup>.

## 4.2 Evaluating Disentanglement

In addition to controllability, disentanglement of the user representations can also be evaluated since the ground truth generative factors are available. Several metrics have been proposed to evaluate disentanglement, such as the Mutual Information Gap (MIG) [8], the  $\beta$ -VAE metric [12] or the FactorVAE metric [15]. We evaluate disentanglement with the Mutual Information Gap (MIG) [8], due to ease of implementation, wide applicability and the unbiased nature of the metric for all hyperparameter settings [20]. Given a generative factor, the empirical mutual information (M.I) between the values of the ground truth generative factors and each dimension of encoded samples from the data is computed. The MIG of this generative factor is then the difference between the highest and second highest MI values. This quantity is averaged across all generative factors to obtain the MIG score. A MIG of 1 for a particular generative factor implies that one latent dimension has  $\text{MI}=1$ , with others having  $\text{MI}=0$ .

## 5 EXPERIMENTS

The datasets, along with the requisite processing steps and generative factors for supervision are described in Section 5.1. The research questions are described in Section 5.2, followed by specifics of models in Section 5.3.

### 5.1 Datasets

In our experiments we use the Million Songs Dataset (MSD)[4] and the Movielens-20M (ML-20M) [10], two widely used collaborative filtering datasets. The steps for preparing the dataset for training are outlined in Section 5.1.1. Preference distributions are constructed as signals for supervision and for evaluating disentanglement/controllability, which is described in Section 5.1.2.

**5.1.1 Dataset processing.** Each dataset is processed according to the steps outlined in [19]. The users are split into test, validation and train sets. The test/validation set size for ML-20M is 10,000 users and for MSD is 50,000 users, with 20% of items held-out. The models is trained with the entire train history, and evaluated with the 20% held out set. The data is binarized by keeping ratings of four or higher and only users who have interacted with at least 5 movies are kept.

**5.1.2 Generative Factors.** The generative factors considered for both datasets are *preference distributions* computed for each user using genres (MSD/ML-20M) or tags (MSD only). The number of users, items, generative factors are reported in Table 2. For a user  $u$ ,  $\text{Genre}_u$  is computed by calculating the proportion of movies that belong to a genre (or tag), divided by the total number of movies that the user has watched, capturing the (global) preferences of a user towards a genre/tag. Since there are 522362 tags in MSD with many repeats, only the most frequent are picked and grouped

<sup>2</sup>Co-occurrences are computed at user level and may differ across users

**Table 2: Datasets used, along with signals being considered. The last column reports the MIG score, computed on the ground truth labels against itself using 10,000 users sampled from the test set**

Dataset (ID)	# users	# items	Generative Factors	G.T MIG
ML-20M	136,677	20108	genre distribution	0.2019
MSD (Genre)	571,355	41140	genre distribution	0.3238
MSD (Tag)	571,355	41140	tag distribution	0.2525

together before computing the distribution<sup>3</sup>. The dimensions of the signals for ML-20M is 19, and for MSD (Genre) is 21, and for MSD(Tag) is 30, with the supervision constraining only the first few dimensions.

## 5.2 Research Questions

This section details the experimental setup employed in the paper. The evaluation metric for recommender performance throughout this paper is NDCG@100, following [19, 30] i.e R-MET is NDCG for all experiments. Controllability for models with 0 supervision cannot be evaluated since  $k(\cdot)$  is unavailable. In addition, the manipulations (which replaces  $[\bar{z}]_{k(q_j)}$ ) are produced by using the Inverse CDF of the probability distribution used during training [17].

**5.2.1 RQ1** *How much supervision is needed for achieving controllability? How does controllability vary across datasets and models?* To investigate this, we train Mult- $\beta$ -VAE and Mult- $\beta$ -VAE with varying levels of supervision:  $\{0\%, 1\%, 10\%, 50\%, 100\%\}$ . For each combination of model, signal, level of supervision, we compute the five metrics outlined earlier. To summarize,  $\delta_{\text{ctrl}}$  evaluates the performance gain when a user sets a knob to its maximum setting, Corr measures the gradual change via a correlation of manipulations and R-MET; the other metrics contrast the controllability of a genre against two control genres. The holdout sets used to compute the metrics above are constructed for 100 users per genre (metrics remain constant beyond 100), while ensuring that there are at least 10/5 items in the input/holdout sets described in Table 1. The number of steps taken in the latent space is 50 for computing the correlation metrics.

**5.2.2 RQ2** *To what extent does disentanglement affect controllability of models?* The relationship between disentanglement and controllability is explored in this section, where we measure if models with disentanglement ( $\beta \geq 1$ ) perform better than models without ( $\beta = 0$ ). Note that  $\beta = 0$  implies no disentanglement, but for the Mult-TCVAE, only the TC constraint is set to 0, i.e which means  $\beta = 0, \gamma = \alpha = 1$ . In addition, we investigate if models that have higher disentanglement (MIG) scores perform better based on recommendation/ controllability metrics. MIG is computed using code from [20]<sup>4</sup>, modified for use in pytorch. Since the exact M.I is intractable, the empirical M.I is computed by discretizing the

values, following [8, 20]. The number of bins for the discretizer is set to 20, and it is computed for on 10,000 samples from the test set.

**5.2.3 RQ3** *What effect does introducing controllability into a model have on recommendation performance?* This question deals with the *relevance* of both controlled and ‘default’ recommendations. This is paramount, otherwise recommendations can be irrelevant. As such, NDCG@100 is computed for all models, and model selection is done on the basis of recommendation performance on the validation set. We remind the readers that  $\delta_{\text{irrel}}$  has to be interpreted with  $\delta_{\text{ctrl}}$ . That is, if  $\delta_{\text{irrel}}$  is low and  $\delta_{\text{ctrl}}$  is high, recommended items may be personalized; however, if  $\delta_{\text{irrel}}$  is higher, the models do exhibit controllability, but the recommended items might not be all personalized. Note that these metrics might be limited since a user might like these items regardless, or, such recommendations might be explicitly sought by the user. Performance for models with no disentanglement ( $\beta = 0$ ) and no supervision ( $\gamma = 0$ ) are also reported.

## 5.3 Experimental Setup

This section details specifics of models being used, along with hyperparameters:

**5.3.1 Variational Distribution.** Both Mult- $\beta$ -VAE and Mult-TCVAE use a isotropic Gaussian distribution for a prior ( $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{1})$ ), and a Gaussian distribution for the variational posterior  $q_\phi$ :

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_\phi(\mathbf{x}), \text{diag}(\boldsymbol{\sigma}_\phi^2(\mathbf{x}))) \quad (7)$$

where  $\boldsymbol{\mu}_\phi(\mathbf{x})$  and  $\boldsymbol{\sigma}_\phi^2(\mathbf{x})$  are outputs of the encoder network.

**5.3.2 Multinomial loss.** We use the Multinomial log-likelihood for  $p_\theta$ , which has been shown to perform well over other losses like the Gaussian/Logistic losses [19]:

$$\log p_\theta(\mathbf{x}|\mathbf{z}) = \sum_i \mathbf{x}_i \log \pi_i(\mathbf{z}) \quad (8)$$

The final loss used for both Mult- $\beta$ -VAE and Mult-TCVAE methods are given in Equation 4, where  $\mathcal{L}_{\text{unsup}}$  is either Equation 2 for the Mult- $\beta$ -VAE or Equation 3 for the Mult-TCVAE, and second term is the binary cross entropy loss (Equation 5). In addition,  $\gamma_{ss} = 0$  for models with no supervision ( $\gamma=0$ ), and  $\gamma_{ss} = 1$  otherwise ( $\gamma>0$ ). To obtain recommendations, items are sorted in the descending order of the likelihood predicted by the decoder.

**5.3.3 Model Hyperparameters.** All models use the same neural network with a 3 layer encoder, with dimensions  $|\mathcal{I}| \rightarrow 600 \rightarrow 600 \rightarrow 200$  and a single layer decoder which was found optimal for a variety of settings [30]. The outputs of the encoder are two 200 dimensional real vectors representing the mean/log-variance vectors. All models were trained with Adam [16] with a learning rate of 0.001 and a batch size of 500 for 200 epochs for ML-20M and 100 epochs for MSD. For both, we tried  $\beta \in \{1, 2.5, 5, 10, 100\}$ , with the ‘best’ model selected using NDCG@100 computed on the validation set. For the Mult-TCVAE, we set  $\alpha = \gamma = 1$ .

<sup>3</sup>Due to a lack of space, these are not included in the paper, and can be accessed at [datasets/msd/selected\\_tags.tsv](https://github.com/google-research/disentanglement_lib/datasets/msd/selected_tags.tsv)

<sup>4</sup>[https://github.com/google-research/disentanglement\\_lib/](https://github.com/google-research/disentanglement_lib/)



**Table 3: Results for ML-20M: MIG measures disentanglement,  $\delta_{\text{ctrl}}$  and Corr measure control, with EasyCorr<sub>ctrl</sub>, EasyCorr<sub>rand</sub>, DiffCorr<sub>ctrl</sub> and DiffCorr<sub>rand</sub> contrasting control of one genre against another. NDCG and  $\delta_{\text{irrel}}$  (when contrasted with  $\delta_{\text{ctrl}}$ ) measure personalization of non- and controlled recommendations respectively. Values in braces are standard errors.**

Model	%	$\beta$	NDCG	MIG	$\uparrow \delta_{\text{ctrl}}$	$\downarrow \delta_{\text{irrel}}$	$\uparrow \text{Corr}$	$\uparrow \text{EasyCorr}_{\text{ctrl}}$	$\downarrow \text{EasyCorr}_{\text{rand}}$	$\uparrow \text{DiffCorr}_{\text{ctrl}}$	$\downarrow \text{DiffCorr}_{\text{rand}}$
Mult- $\beta$ -VAE	0	0	0.2782	0.0016	-	-	-	-	-	-	-
	50	0	0.2627	0.1265	0.1995 (0.0153)	0.0658 (0.0100)	0.8859	0.8880 (0.0045)	-0.2487 (0.0132)	0.8547 (0.0055)	-0.5793 (0.0115)
	100	0	0.2587	0.1417	<b>0.2246</b> (0.0190)	0.0764 (0.0127)	0.8628	0.8677 (0.0053)	-0.1994 (0.0119)	0.8408 (0.006)	-0.5911 (0.0113)
Mult-TCVAE	0	0	<b>0.2904</b>	0.0007	-	-	-	-	-	-	-
	50	0	0.2802	0.0688	0.2145 (0.0157)	0.0705 (0.0104)	<b>0.8901</b>	0.8907 (0.0044)	-0.2869 (0.0125)	0.8599 (0.0053)	-0.6448 (0.0104)
	100	0	0.2791	<b>0.1819</b>	0.2210 (0.0198)	0.0755 (0.0129)	0.8522	0.8557 (0.0055)	-0.2054 (0.0113)	0.8357 (0.0061)	-0.6569 (0.0101)
Mult- $\beta$ -VAE	0	2.5	0.3740	0.0184	-	-	-	-	-	-	-
	1	100	0.3807	0.0208	0.0904 (0.0091)	0.0313 (0.0049)	0.8457	0.8449 (0.0059)	-0.2010 (0.0160)	0.7647 (0.0083)	-0.5283 (0.0125)
	10	2.5	0.3808	0.0858	0.1863 (0.0128)	0.0644 (0.0091)	<b>0.9016</b>	0.9039 (0.0038)	-0.2664 (0.0134)	0.8567 (0.0053)	-0.6208 (0.0109)
	50	2.5	0.3703	0.1360	0.2314 (0.0173)	0.0800 (0.0122)	0.8886	0.8893 (0.0046)	-0.2951 (0.0122)	0.8506 (0.0057)	-0.6414 (0.0104)
	100	10	0.3675	0.1476	0.2215 (0.0188)	0.0770 (0.0132)	0.8632	0.8681 (0.0053)	-0.2304 (0.0119)	0.8337 (0.0062)	-0.6483 (0.0106)
Mult-TCVAE	0	1	0.3897	0.0062	-	-	-	-	-	-	-
	1	1	<b>0.3898</b>	0.0071	-0.0123 (0.0187)	0.0321 (0.0092)	0.2564	0.2791 (0.0146)	-0.1947 (0.0140)	0.3141 (0.0143)	-0.1164 (0.0159)
	10	1	0.3884	0.0087	0.1283 (0.0222)	0.0991 (0.0354)	0.8662	0.8735 (0.0035)	-0.3599 (0.0131)	0.8542 (0.0042)	-0.2949 (0.0142)
	50	1	0.3876	0.0597	0.2885 (0.0327)	0.2663 (0.0610)	0.8584	0.8648 (0.0053)	-0.3617 (0.0125)	0.9007 (0.0029)	-0.5658 (0.0138)
	100	1	0.3863	<b>0.1967</b>	<b>0.3844</b> (0.0381)	0.3629 (0.0657)	0.8279	0.8379 (0.006)	-0.2221 (0.0110)	0.8806 (0.0039)	-0.5880 (0.0143)

**Table 4: Results for MSD (Genre): Models for this dataset in particular struggle to distinguish between genres (very high CorrCtrl scores)**

Model	%	$\beta$	$\uparrow \text{NDCG}$	MIG	$\uparrow \delta_{\text{ctrl}}$	$\downarrow \delta_{\text{irrel}}$	$\uparrow \text{Corr}$	$\uparrow \text{EasyCorr}_{\text{ctrl}}$	$\downarrow \text{EasyCorr}_{\text{rand}}$	$\uparrow \text{DiffCorr}_{\text{ctrl}}$	$\downarrow \text{DiffCorr}_{\text{rand}}$
Mult- $\beta$ -VAE	0	0	0.2537	<b>0.0027</b>	-	-	-	-	-	-	-
	50	0	0.2503	0.0022	0.0974 (0.0068)	0.0680 (0.0119)	0.6106	0.7150 (0.0119)	0.2413 (0.0212)	0.7251 (0.0117)	0.2745 (0.0235)
	100	0	0.2459	0.0022	0.1378 (0.011)	0.0972 (0.0171)	0.6191	0.7173 (0.0115)	0.2448 (0.0205)	0.7462 (0.0110)	0.2836 (0.0233)
Mult-TCVAE	0	0	<b>0.2620</b>	0.0022	-	-	-	-	-	-	-
	50	0	0.2543	0.0022	0.0979 (0.0095)	0.0715 (0.0120)	0.6021	0.7122 (0.0120)	0.2309 (0.0213)	0.7334 (0.0116)	0.2688 (0.0240)
	100	0	0.2509	0.0022	<b>0.1614</b> (0.0118)	0.1087 (0.0173)	<b>0.6522</b>	0.7492 (0.0108)	0.2612 (0.0211)	0.7715 (0.0103)	0.2674 (0.0243)
Mult- $\beta$ -VAE	0	5	0.2768	0.0055	-	-	-	-	-	-	-
	1	100	0.2759	0.0076	0.0827 (0.0242)	0.1535 (0.0331)	0.5111	0.5659 (0.0156)	0.2417 (0.0207)	0.5739 (0.0159)	0.0855 (0.0232)
	10	1	0.2745	0.0194	0.3561 (0.0527)	0.4560 (0.0177)	0.7034	0.6824 (0.0118)	0.2903 (0.0201)	0.6907 (0.0117)	0.1200 (0.0231)
	50	10	0.2706	0.0388	0.4693 (0.0382)	0.6822 (0.0202)	<b>0.7807</b>	0.8014 (0.0048)	0.2587 (0.0218)	0.8003 (0.0050)	0.1307 (0.0255)
	100	1	0.2719	<b>0.4178</b>	<b>0.4853</b> (0.0342)	0.7795 (0.0322)	0.7209	0.7596 (0.0062)	0.1894 (0.0222)	0.7598 (0.0062)	0.0613 (0.0266)
Mult-TCVAE	0	1	0.2763	0.0145	-	-	-	-	-	-	-
	1	1	<b>0.2777</b>	0.0069	0.0131 (0.0061)	0.0216 (0.0085)	0.2302	0.3328 (0.0195)	0.1143 (0.0200)	0.3352 (0.0202)	0.1461 (0.0220)
	10	1	0.2728	0.0148	0.3461 (0.044)	0.3915 (0.0296)	0.7343	0.7501 (0.0082)	0.2040 (0.0226)	0.7573 (0.0082)	0.1680 (0.0241)
	50	1	0.2675	0.0471	0.4439 (0.0366)	0.5324 (0.0419)	0.7820	0.8001 (0.0059)	0.2405 (0.0224)	0.8099 (0.0057)	0.1438 (0.0259)
	100	1	0.2671	0.3521	0.4501 (0.0328)	0.7526 (0.0405)	0.6772	0.7219 (0.0073)	0.1927 (0.021)	0.7057 (0.0080)	0.0539 (0.0257)

## 6 RESULTS

The results of all experiments on ML-20M is reported in Table 3, on MSD (Tag) in Table 4 and finally on MSD (Genre) in Table 5. The results indicate that while increasing supervision generally tends to increase controllability, too much supervision can sometimes harm controllability, or alternatively increase controllability at the cost of more non-personalized recommendations (Section 6.1). The degree of supervision required for controllability varies across models/datasets/signals. In addition, adding a disentanglement objective helps across all metrics, but there appears to be no conclusive trend between degree of disentanglement and controllability (Section 6.2). Finally, controlled recommendations seem to be more personalized for ML-20M compared to MSD, where more non-personalized items

get recommended (Section 6.3). We discuss each result in detail in the sections that follow.

### 6.1 Semi-Supervision and Controllability

Surprisingly, apparent control over recommendations seems to be achieved with as *little as 10% of data* for both datasets, as measured by high Corr and  $\delta_{\text{ctrl}}$ . In addition, these values increase as supervision is increased to 50% for all data and models. Supervision beyond that, however, can produce mixed results; for instance, Mult- $\beta$ -VAE trained on 100% of ML-20M scores worse on  $\delta_{\text{ctrl}}$  and Corr. In contrast, for MSD(Genre) and MSD(Tag) (with the exception of Mult- $\beta$ -VAE for 100%),  $\delta_{\text{ctrl}}$  increases with supervision. These trends also hold for the  $\beta = 0$  models. Therefore, *in most cases, supervision increases apparent controllability*.

**Table 5: Results for MSD (Tag): Models confuse between tags to a lesser extent compared to Genre**

Model	%	$\beta$	$\uparrow$ NDCG	MIG	$\uparrow \delta_{\text{ctrl}}$	$\downarrow \delta_{\text{irrel}}$	$\uparrow$ Corr	$\uparrow$ EasyCorr <sub>ctrl</sub>	$\downarrow$ EasyCorr <sub>rand</sub>	$\uparrow$ DiffCorr <sub>ctrl</sub>	$\downarrow$ DiffCorr <sub>rand</sub>
Mult- $\beta$ -VAE	0	0	0.2537	0.0006	-	-	-	-	-	-	-
	50	0	0.2419	0.0006	0.0832 (0.0035)	0.0835 (0.0092)	<b>0.7608</b>	0.7622 (0.0056)	-0.0824 (0.0068)	0.7168 (0.0061)	-0.2136 (0.0118)
	100	0	0.2370	0.0006	0.0776 (0.0050)	0.0790 (0.0111)	0.6812	0.6804 (0.0064)	-0.0405 (0.0052)	0.6685 (0.0065)	-0.2182 (0.0118)
Mult-TCVAE	0	0	<b>0.2620</b>	0.0006	-	-	-	-	-	-	-
	50	0	0.2444	0.0006	<b>0.0905</b> (0.0039)	0.0925 (0.0097)	0.7546	0.7577 (0.0055)	-0.0638 (0.0058)	0.7292 (0.0058)	-0.1664 (0.0123)
	100	0	0.2403	0.0006	0.0828 (0.0044)	0.0807 (0.0100)	0.6960	0.6989 (0.0062)	-0.0539 (0.0053)	0.6858 (0.0062)	-0.1941 (0.0121)
Mult- $\beta$ -VAE	0	2.5	0.2768	0.0038	-	-	-	-	-	-	-
	1	1	0.2749	0.0107	0.0407 (0.0095)	0.2519 (0.0358)	0.5815	0.5798 (0.0083)	-0.0865 (0.0083)	0.5827 (0.0082)	-0.0230 (0.0115)
	10	1	0.2737	0.0233	0.1572 (0.0120)	0.5032 (0.0310)	0.7946	0.7954 (0.0041)	-0.1542 (0.0083)	0.8074 (0.0039)	-0.1773 (0.0123)
	50	1	0.2682	0.0847	0.1886 (0.0143)	0.5553 (0.0455)	0.8194	0.8161 (0.0039)	-0.1194 (0.0075)	0.8141 (0.0038)	-0.1912 (0.0126)
	100	5	0.2655	0.2244	0.1874 (0.0121)	0.5763 (0.0473)	0.8170	0.8192 (0.0037)	-0.1010 (0.007)	0.8210 (0.0038)	-0.1283 (0.0130)
Mult-TCVAE	0	1	0.2781	0.0072	-	-	-	-	-	-	-
	1	1	<b>0.2822</b>	0.0085	0.0165 (0.0042)	0.0741 (0.0160)	0.4415	0.4501 (0.0095)	-0.0359 (0.0088)	0.4498 (0.0095)	-0.0042 (0.0112)
	10	1	0.2761	0.0065	0.0925 (0.0086)	0.2490 (0.0404)	0.7915	0.7948 (0.0048)	-0.1374 (0.0086)	0.7872 (0.0049)	-0.0929 (0.0118)
	50	1	0.2670	0.1123	0.2018 (0.0151)	0.4289 (0.0437)	<b>0.8481</b>	0.8494 (0.0036)	-0.1290 (0.0084)	0.8431 (0.0036)	-0.1133 (0.0130)
	100	1	0.2628	<b>0.2272</b>	<b>0.2029</b> (0.0133)	0.5042 (0.0454)	0.8224	0.8224 (0.0041)	-0.0709 (0.0062)	0.8237 (0.0039)	-0.1316 (0.0130)

Given a model which scores high on  $\delta_{\text{ctrl}}$  (or Corr), finer grained RandCorr metrics can be analyzed, to check if changing a knob inadvertently changes other genres<sup>5</sup>. Increasing supervision for ML-20M makes models score better on DiffCorr<sub>rand</sub> (i.e more negative), while only slightly decreasing DiffCorr<sub>ctrl</sub>, indicating that supervision can help with similar genres in particular. While supervision initially improves EasyCorr<sub>rand</sub>, 100% supervision actually worsens it slightly from 50%. Therefore, for ML-20M, supervision can help the model avoid confusing similar genres, at the (slight) cost of controlling dissimilar genres, especially if fully supervised.

For MSD (Genre), we first note that values of EasyCorr<sub>rand</sub> and DiffCorr<sub>rand</sub> are positive for this dataset, indicating that other genres are also being manipulated. Increasing supervision increases EasyCorr<sub>ctrl</sub> and DiffCorr<sub>ctrl</sub> i.e help control, while EasyCorr<sub>rand</sub> and DiffCorr<sub>rand</sub> tend to become worse, indicating confusion among other genres. While EasyCorr<sub>ctrl</sub> and DiffCorr<sub>ctrl</sub> becomes worse with 100% supervision, EasyCorr<sub>rand</sub> and DiffCorr<sub>rand</sub> improve, indicating that full supervision can help with confusion. In conclusion, for MSD (Genre), increasing supervision generally improves control, at the cost of other genres being also recommended, which improves with 100% supervision for MSD(Genre).

In MSD (Tag), supervision generally increases EasyCorr<sub>ctrl</sub> and DiffCorr<sub>ctrl</sub>. For the Mult-TCVAE, increasing supervision to 100% makes the model score better on DiffCorr<sub>rand</sub>, but worse on EasyCorr<sub>rand</sub>. In contrast, for increasing supervision to 100% for the Mult- $\beta$ -VAE makes the model score worse on both DiffCorr<sub>rand</sub> and EasyCorr<sub>rand</sub>. For the MSD (Tag) dataset, therefore, while supervision increases controllability, too much can cause confusion between tags.

In summary, *while an increase in supervision generally increases controllability across datasets and models, and additionally reduces confusion between factors, too much supervision can harm performance* by causing other genres to also be recommended alongside the one being controlled.

<sup>5</sup>Note that EasyCorr<sub>ctrl</sub> (or EasyCorr<sub>rand</sub>) cannot be compared with DiffCorr<sub>ctrl</sub> (or DiffCorr<sub>rand</sub>), as the inputs for computing these metrics can be very different. In addition, we observed that a lot of values of EasyCorr<sub>rand</sub> (at user/genre level) are equal to 0.

## 6.2 Disentanglement and Controllability

This section discusses if disentanglement ( $\beta \geq 1$ ) models are necessary, and if increased disentanglement as measured by the MIG score contributes to controllability.

**6.2.1 Comparison with baselines.** Models without disentanglement  $\beta = 0$ , generally score worse in most respects, compared to models with disentanglement. For the ML-20M dataset,  $\beta = 0$  models have comparable controllability scores compared with models with disentanglement. However, this comes at a drastic reduction in recommender performance i.e the best NDCG (0.2904) among the baselines is worse than the lowest NDCG among models with disentanglement (0.3675). This is also seen for the MSD (Genre) datasets, although the performance drop is smaller. For both the MSD datasets, however, the gap observed in the controllability performance of entangled/disentangled models is much greater than the gap in the ML-20M dataset, indicating that disentanglement is necessary for controllability in these datasets. Overall, *models with disentanglement perform better than models without disentanglement*. The next section compares the *degree* of disentanglement with controllability.

**6.2.2 Does increased disentanglement help?** For the ML-20M dataset, increased disentanglement results in lower recommendation (NDCG) performance. ML-20M models with high MIG scores perform better on  $\delta_{\text{ctrl}}$ , although this doesn't seem to be necessary for Mult-TCVAE with  $\beta = 50$ , which scores high on many metrics despite having a low MIG score. For the MSD (Genre) models, a jump in MIG is accompanied by large increases in  $\delta_{\text{ctrl}}$  and  $\delta_{\text{irrel}}$ . However, as the values of MIG scores here are either high or near zero, it's difficult to make concrete conclusions about the precise relationship between disentanglement scores and controllability. We note, however, that Mult-TCVAE, which achieves better disentanglement [8] in general, seems to perform better on the ML-20M and MSD (Tag) dataset while achieving the highest MIG scores, while Mult- $\beta$ -VAE performs better and achieves the highest MIG score for MSG (Genre).



### 6.3 Controlled Recommendation Performance

We first note that *as the amount of supervision, and consequently, controllability increases, average recommender performance (NDCG) decreases*. However, we argue that the benefits might outweigh the cost, as this drop is relatively small while providing control for users (or practitioners). While NDCG measures performance of non-controlled recommendations, the degree of personalization of controlled recommendations can be measured by comparing  $\delta_{\text{ctrl}}$  and  $\delta_{\text{irrel}}$ .

The results vary depending on the model or dataset. For instance, since  $\delta_{\text{ctrl}}$  is much higher than  $\delta_{\text{irrel}}$  for Mult- $\beta$ -VAE with 100% supervision on the ML-20M, we can conclude that it recommends personalized items instead of random items that a user might like. However, the opposite is true for Mult-TCVAE, as  $\delta_{\text{ctrl}}$  and  $\delta_{\text{irrel}}$  are at a similar level, which means that the items recommended on controlling a genre might not be as personalized to the user. Mult-TCVAE produces the highest  $\delta_{\text{ctrl}}$  scores for the ML-20M and MSD(Tag) datasets, while Mult- $\beta$ -VAE produces the highest  $\delta_{\text{ctrl}}$  scores for MSD (Genre). In addition, while  $\delta_{\text{ctrl}}$  is higher (more items of that genre are being pushed to the top),  $\delta_{\text{irrel}}$  is also very high in MSD (Genre), which means recommendations are controllable, but might not be personalized. However, as data sparsity increases,  $\delta_{\text{irrel}}$  might be overestimated because of missing relevance judgements, as evidenced by high  $\delta_{\text{irrel}}$  for almost all models trained on MSD. In addition to incompleteness, this evaluation also does not account for the exploratory/interactive/serendipitous nature of controllable recommendations.

In summary, *controlled recommendations appear to be personalized*, as evidenced by high scores of  $\delta_{\text{ctrl}}$ , showing that users can exert control and expect personalized recommendations to an extent. However, *other items of the given genre might also be recommended*, which is especially true for the MSD dataset.

## 7 CONCLUSION

We showed that controllable recommendations can be achieved by leveraging the generative nature of recommenders. (Semi-)supervised disentanglement is used both to tie a generative factor to a known dimension, and to enforce disentanglement. Consequently, manipulating a dimension produces *controlled* recommendations, allowing a user to express short-term/dynamic preferences, or to explore these recommendations in an interactive manner. We experimented with genre/tag distributions as supervision signals on two datasets. We proposed metrics to measure the extent of controllability in addition to the degree of personalization of controlled recommendations. Using these metrics, we showed that a user can control recommendations, while retaining personalization to an extent. We showed that such control comes only at a slight reduction in recommendation performance, but might require different degrees of supervision depending on the data, model or signals used.

While the experiments here detail controlling only a single dimension, multiple dimensions can be manipulated allowing for greater or more nuanced control, which we leave for future work. In addition, the dimensionality of the supervision signal is limited by the dimensionality of the latent space, which limits the applicability to large sets of knobs. Evaluating the personalization of controlled recommendations presents a challenge as outlined in

the previous section, which we plan to pursue in the future. In addition, this method can be used to ‘boost’ recommendations of items belonging to a particular category, for instance, if items of category  $X$  are being under-recommended (possibly due to biases), they can be boosted to compensate for it.

## ACKNOWLEDGMENTS

The authors would like to thank Mohammad Aliannejadi, Wilker Aziz, Maurits Bleeker, Jin Huang, Antonis Krasakis, Anna Sepiarskaia, Georgios Sidiropoulos and Svitlana Vakulenko for helpful comments and feedback. The authors also thank the feedback received by reviewers during prior review processes. This research was supported by the NWO Innovational Research Incentives Scheme Vidi (016.Vidi.189.039), the NWO Smart Culture - Big Data / Digital Humanities (314-99-301), and the H2020-EU.3.4. - SOCIETAL CHALLENGES - Smart, Green And Integrated Transport (814961). All content represents the opinion of the authors, which is not necessarily shared or endorsed by their respective employers and/or sponsors.

## REFERENCES

- [1] Diego Antognini and Boi Faltings. 2021. Fast Multi-Step Critiquing for VAE-Based Recommender Systems. In *Fifteenth ACM Conference on Recommender Systems* (Amsterdam, Netherlands) (*RecSys '21*). Association for Computing Machinery, New York, NY, USA, 209–219. <https://doi.org/10.1145/3460231.3474249>
- [2] Diego Antognini, Claudiu Musat, and Boi Faltings. 2020. T-RECS: a Transformer-based Recommender Generating Textual Explanations and Integrating Unsupervised Language-based Critiquing. *CoRR* abs/2005.11067 (2020). arXiv:2005.11067 <https://arxiv.org/abs/2005.11067>
- [3] Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2013. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence* 35, 8 (2013), 1798–1828.
- [4] Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. 2011. The Million Song Dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*.
- [5] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. 2017. Variational inference: A review for statisticians. *Journal of the American statistical Association* 112, 518 (2017), 859–877.
- [6] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. 2018. Understanding disentangling in  $\beta$ -VAE. *arXiv preprint arXiv:1804.03599* (2018).
- [7] Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. Controllable Multi-Interest Framework for Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2942–2951.
- [8] Ricky TQ Chen, Xuechen Li, Roger B Grosse, and David K Duvenaud. 2018. Isolating sources of disentanglement in variational autoencoders. In *Advances in Neural Information Processing Systems*. 2610–2620.
- [9] Zeyu Cui, Feng Yu, Shu Wu, Qiang Liu, and Liang Wang. 2020. Disentangled Item Representation for Recommender Systems. *arXiv preprint arXiv:2008.07178* (2020).
- [10] F. Maxwell Harper and Joseph A. Konstan. 2015. The MovieLens Datasets: History and Context. *ACM Trans. Interact. Intell. Syst.* 5, 4, Article 19 (Dec. 2015), 19 pages. <https://doi.org/10.1145/2827872>
- [11] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.
- [12] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. 2017.  $\beta$ -VAE: Learning Basic Visual Concepts with a Constrained Variational Framework. *Iclr* 2, 5 (2017), 6.
- [13] Dietmar Jannach, Ahtsham Manzoor, Wanling Cai, and Li Chen. 2020. A survey on conversational recommender systems. *arXiv preprint arXiv:2004.00646* (2020).
- [14] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. 1999. An introduction to variational methods for graphical models. *Machine learning* 37, 2 (1999), 183–233.
- [15] Hyunjik Kim and Andriy Mnih. 2018. Disentangling by factorising. *arXiv preprint arXiv:1802.05983* (2018).
- [16] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [17] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [18] Hanze Li, Scott Sanner, Kai Luo, and Ga Wu. 2020. A Ranking Optimization Approach to Latent Linear Critiquing for Conversational Recommender Systems. In *Fourteenth ACM Conference on Recommender Systems* (Virtual Event, Brazil) (*RecSys '20*). Association for Computing Machinery, New York, NY, USA, 13–22. <https://doi.org/10.1145/3383313.3412240>
- [19] Dawen Liang, Rahul G Krishnan, Matthew D Hoffman, and Tony Jebara. 2018. Variational autoencoders for collaborative filtering. In *Proceedings of the 2018 World Wide Web Conference*. 689–698.
- [20] Francesco Locatello, Stefan Bauer, Mario Lucic, Gunnar Raetsch, Sylvain Gelly, Bernhard Schölkopf, and Olivier Bachem. 2019. Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations (*Proceedings of Machine Learning Research*, Vol. 97), Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), PMLR, Long Beach, California, USA, 4114–4124. <http://proceedings.mlr.press/v97/locatello19a.html>
- [21] Francesco Locatello, Michael Tschannen, Stefan Bauer, Gunnar Rätsch, Bernhard Schölkopf, and Olivier Bachem. 2019. Disentangling factors of variation using few labels. *arXiv preprint arXiv:1905.01258* (2019).
- [22] Kai Luo, Scott Sanner, Ga Wu, Hanze Li, and Hojin Yang. 2020. Latent Linear Critiquing for Conversational Recommender Systems. In *Proceedings of The Web Conference 2020* (Taipei, Taiwan) (*WWW '20*). Association for Computing Machinery, New York, NY, USA, 2535–2541. <https://doi.org/10.1145/3366423.3380003>
- [23] Kai Luo, Hojin Yang, Ga Wu, and Scott Sanner. 2020. *Deep Critiquing for VAE-Based Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 1269–1278. <https://doi.org/10.1145/3397271.3401091>
- [24] Jianxin Ma, Peng Cui, Kun Kuang, Xin Wang, and Wenwu Zhu. 2019. Disentangled graph convolutional networks. In *International Conference on Machine Learning*. 4212–4221.
- [25] Jianxin Ma, Chang Zhou, Peng Cui, Hongxia Yang, and Wenwu Zhu. 2019. Learning disentangled representations for recommendation. In *Advances in Neural Information Processing Systems*. 5712–5723.
- [26] Jianxin Ma, Chang Zhou, Hongxia Yang, Peng Cui, Xin Wang, and Wenwu Zhu. 2020. Disentangled Self-Supervision in Sequential Recommenders. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 483–491.
- [27] Loic Matthey, Irina Higgins, Demis Hassabis, and Alexander Lerchner. 2017. dSprites: Disentanglement testing Sprites dataset. <https://github.com/deepmind/dsprites-dataset/>.
- [28] Preksha Nema, Alexandros Karatzoglou, and Filip Radlinski. 2021. Disentangling Preference Representations for Recommendation Critiquing with  $\beta$ -VAE. In *30th ACM International Conference on Information and Knowledge Management (CIKM 2021)*. New York.
- [29] Francesco Ricci, Lior Rokach, Bracha Shapira, and Paul B. Kantor. 2010. *Recommender Systems Handbook* (1st ed.). Springer-Verlag, Berlin, Heidelberg.
- [30] Ilya Shenbin, Anton Alekseev, Elena Tutubalina, Valentin Malykh, and Sergey I Nikolenko. 2020. RecVAE: A New Variational Autoencoder for Top-N Recommendations with Implicit Feedback. In *Proceedings of the 13th International Conference on Web Search and Data Mining*. 528–536.
- [31] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. 2015. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 1235–1244.
- [32] Haonan Wang, Chang Zhou, Carl Yang, Hongxia Yang, and Jingrui He. 2021. Controllable Gradient Item Retrieval. In *Proceedings of the Web Conference 2021* (Ljubljana, Slovenia) (*WWW '21*). Association for Computing Machinery, New York, NY, USA, 768–777. <https://doi.org/10.1145/3442381.3449963>
- [33] Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. 2019. Neural graph collaborative filtering. In *Proceedings of the 42nd international ACM SIGIR conference on Research and development in Information Retrieval*. 165–174.
- [34] Xiang Wang, Hongye Jin, An Zhang, Xiangnan He, Tong Xu, and Tat-Seng Chua. 2020. Disentangled Graph Collaborative Filtering. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1001–1010.
- [35] Ga Wu, Kai Luo, Scott Sanner, and Harold Soh. 2019. Deep Language-Based Critiquing for Recommender Systems. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (*RecSys '19*). Association for Computing Machinery, New York, NY, USA, 137–145. <https://doi.org/10.1145/3298689.3347009>
- [36] Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. 2016. Collaborative denoising auto-encoders for top-n recommender systems. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*. 153–162.
- [37] Hojin Yang, Tianshu Shen, and Scott Sanner. 2021. *Bayesian Critiquing with Keyphrase Activation Vectors for VAE-Based Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 2111–2115. <https://doi.org/10.1145/3404835.3463108>
- [38] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM Computing Surveys (CSUR)* 52, 1 (2019), 1–38.