Structure-Conditional Minimum Bayes Risk Decoding

Anonymous ACL submission

Abstract

Minimum Bayes Risk (MBR) decoding has seen renewed interest as an alternative to traditional generation strategies. While MBR has proven effective in machine translation, where the variability of a language model's outcome space is naturally constrained, it may face challenges in more open-ended tasks such as dialogue or instruction-following. We hypothesise that in such settings, applying MBR with standard similarity-based utility functions may result in selecting responses that are broadly rep-011 resentative of the model's distribution, yet suboptimal with respect to any particular grouping of generations that share an underlying la-014 tent structure. In this work, we introduce three 015 lightweight adaptations to the utility, designed 017 to make MBR more sensitive to structural variability in the outcome space. To test our hypothesis, we curate a dataset capturing three representative types of latent structure: dialogue act, emotion, and response structure. We also propose two metrics to evaluate the structural optimality of MBR. Our analysis demonstrates that common utility functions fall short by these metrics. In contrast, our proposed adaptations considerably improve structural optimality. Finally, we evaluate our approaches on real-world 027 instruction-following benchmarks, AlpacaEval and MT-Bench, and show that increased structural sensitivity improves generation quality by up to 13.7 percentage points in win rate.

1 Introduction

Once a language model has been trained, one fundamental problem remains: determining how to select an output sequence from the model's learned probability distribution over possible continuations. Traditional approaches such as beam search decoding and majority voting aim to select a high probability continuation under the model distribution. However, a growing body of research has shown model probability to not reliably align with human preferences (Stahlberg and Byrne, 2019; Zhang



Figure 1: The choice of utility function can considerably impact the Minimum Bayes Risk (MBR) optimum. When the outcome space is structured or multimodal, the MBR optimum may settle between modes, landing in a region of low probability. Here, we present a continuous example featuring a bimodal Gaussian distribution and show the MBR optima (vertical lines) of two utility functions with markedly different behaviours.

043

044

045

046

047

050

051

054

059

061

062

063

064

065

et al., 2021) and, in response, Minimum Bayes Risk (MBR; Kumar and Byrne, 2004; Eikema and Aziz, 2020) decoding has emerged as a more robust alternative. MBR casts decoding as a decision-theoretic problem, selecting the output that minimises risk with respect to a task-specific utility function, under the uncertainty over continuations represented by the language model. This utility typically reflects the degree of agreement between a candidate and the broader set of outcomes, penalising those that diverge significantly from the consensus. By integrating both model probabilities and intercandidate consistency, MBR yields generations that are better aligned with human preferences, regularly outperforming conventional methods (Freitag et al., 2022; Wu et al., 2025).

MBR decoding has gained significant attention in neural machine translation, where utility is often measured by task-agnostic sentence similarity scores. This corresponds to selecting the sequence which, in expectation and under the lens of a particular similarity score, most closely matches the broader distribution of sequences prescribed by the candidate's (soft) membership in structure-specific candidate groups, while preserving the decisiontheoretic foundation of risk minimisation. Our experiments confirm that adapting the utility function to account for latent structural variability substantially improves MBR solutions. On our curated dataset with controlled uncertainty over dialogue act, emotion, and response structure, our three proposed methods achieve markedly higher cluster optimality than standard MBR with BERTScore or BLEURT utilities. We also observe gains on real-world instruction-following benchmarks, despite the absence of explicit structure annotations. In particular, our methods improve generation quality on AlpacaEval and MT-Bench, with win rates against GPT-40 increasing by up to 13.7 percentage points on the latter. These findings support our central claim: structure-aware utility functions enable MBR to more reliably select highquality outputs in tasks where structural variability is inherent to the generation space.

Clustering, Structure Embeddings, and Utility Cut-

off-that adapt utility functions to account for a

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

2 Language Modeling and Decision Rules

A language model P is a distribution over strings Σ^* , where Σ is an alphabet, *i.e.*, a finite, non-empty set of symbols, and Σ^* its Kleene closure, *i.e.*, the set of all strings formed by concatenating symbols in Σ , including the empty string ε . We define Y as a random variable over sequences within Σ^* . Every language model can be expressed in autoregressive form by decomposing the probability of a string as the product of conditional probabilities of each symbol, followed by an end-of-string event EOS:

$$P(Y = y) = P(\text{EOS} \mid y) \prod_{t=1}^{|y|} P(y_t \mid y_{< t}), \quad (1)$$

where each conditional distribution $P(Y_t \mid y_{\le t})$ 152 is a probability distribution over $\Sigma \cup \{EOS\}$. This 153 formulation underlies most modern autoregressive 154 language models, where each conditional probabil-155 ity is produced by a learned parametric model. We 156 will assume an implicit conditioning on some set 157 of neural network parameters θ estimated during 158 neural network training on some dataset. Further-159 more, a language model is commonly conditioned 160 on an input, or a prompt, x. We are then interested 161 only in the conditional probability distribution over responses P(Y|x). In the rest of this work, we 163

model. While this decoding strategy works well 066 in translation-where outcome space variability 067 is inherently constrained by the task-it risks be-068 ing less effective for tasks with a broader range of contextually plausible latent structures, and thus greater variability in realisations, such as dialogue 071 or instruction-following (Giulianelli et al., 2023). 072 Consider, for instance, the following dialogue exchange: A: The mountains would be a great place for the lab retreat. B: That's a wonderful choice. In response, speaker A could follow up with a statement (The mountains offer many outdoor team-077 building activities.), a question (Which aspects of the mountains are you most excited about?), an instruction (Please check out different venues online to finalise the decision.), or an offer (Shall I make the necessary arrangements?). Similarly, when given an instruction like Please summarise Gödel, Escher, Bach, valid responses could range from a single-sentence summary, to a paragraph with more detail, a multi-paragraph narrative, or even a list of key topics. In such settings, applying MBR with a standard similarity-based utility function may result in selecting an output that is broadly representative of the model's outcome distribution, but suboptimal with respect to any one plausible 091 latent structure (we illustrate this in Fig. 1 using a continuous distribution as a simplified example).

In this work, we propose adapting utility functions for MBR such that they are able to explicitly account for a language model's uncertainty over 096 latent structures. We adopt a broad interpretation 098 of structure, treating it as a latent variable that influences the form a generation takes—for example, a dialogue act, the level of detail in a response, or 100 the emotion conveyed by an utterance. To examine 101 how reliably MBR selects the highest-consistency 102 candidate within clusters of generations that share 103 104 a latent structure—what we call *cluster-optimality*, we semi-automatically construct a dataset of 3,000 105 curated outcome spaces, for a total of 350,000 106 candidate generations. These are conditioned on naturally occurring conversational and instruction-108 following contexts, but present controlled uncer-109 tainty over three types of structure: dialogue act, 110 emotion, and response structure (i.e., a single sen-111 tence, a paragraph, a list, or a table). Our anal-112 ysis of this dataset shows that, under commonly 113 used utility functions, MBR solutions are cluster-114 optimal in fewer than half of the cases. To ad-115 dress this, we introduce three new approaches-116

166

176

177

178

179

181

183

184

187

188

189

191

192

193

194

195

197

198

199

200

202

207

will always assume the presence of such an input x, e.g., an instruction or dialogue history.

Decision Rules 2.1

In order to obtain a generation from a trained lan-167 168 guage model P given some input x, it is necessary to decide on a single "best" outcome in Σ^* . 169 Formally, this requires a decision rule that defines 170 a mapping from a distribution P to such an outcome y^* . A common choice is to output the highest 172 probability outcome under P(Y|x), a decision rule 173 known as maximum-a-posteriori, typically approx-174 imated using beam search or majority voting. 175

$$y^*_{\mathsf{MAP}} = \operatorname*{argmax}_{h \in \Sigma^*} P(Y = h | x) \tag{2}$$

However, studies have shown that model probability does not reliably align with human preferences (Stahlberg and Byrne, 2019; Zhang et al., 2021), and Minimum Bayes Risk (MBR) has become a popular alternative. MBR stems from the principle of maximisation of expected utility (Berger, 1985). It requires choosing a *utility* function u(h, r)that measures the benefit of choosing hypothesis h given an ideal decision r. In natural language generation, u is typically chosen to be a strong sentence similarity metric such as BLEURT (Sellam et al., 2020; Freitag et al., 2022), COMET (Rei et al., 2020; Fernandes et al., 2022) or BERTScore (Zhang et al., 2020; Suzgun et al., 2023). MBR then selects the outcome maximising utility in expectation under the model distribution:

$$y_{\text{MBR}}^* = \underset{h \in \Sigma^*}{\operatorname{argmax}} \underset{P(Y|x)}{\mathbb{E}} [u(h, Y)]$$
(3)

Recently, a sampling-based approximation to MBR has become popular that approximates expected utility by obtaining a set of unbiased samples from the model and re-ranking them using Monte Carlo estimates of their expected utility (Eikema and Aziz, 2020, 2022). In this work, we will focus on this sampling-based approximation.

2.2 Structural Variation in Language Models

The importance of modelling uncertainty in natural language generation systems has received growing attention in recent years (Baan et al., 2023). Crucially, uncertainty extends beyond surface-form variations in outcome space to encompass deeper 206 variation in latent space. To capture such variation, metrics like semantic entropy (Kuhn et al., 2024) and similarity-sensitive entropy (Cheng and

Vlachos, 2024) have been proposed, primarily to identify when high uncertainty may signal potential model errors. Complementary work has examined similar measures with a different aim: to assess whether the uncertainty exhibited by language models aligns with the natural variability found in human-generated responses (Deng et al., 2022; Giulianelli et al., 2023; Ilia and Aziz, 2024).

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

Recent applications of MBR have largely focused on neural machine translation-a relatively constrained task where, nonetheless, models have been shown to capture less variation than what human annotators consider plausible (Giulianelli et al., 2023). Extending beyond translation, a few studies have applied MBR to other generation tasks. For example, Suzgun et al. (2023) successfully use BERTScore-based MBR for summarisation, datato-text generation, textual style transfer, and image captioning. However, these tasks also tend to involve a limited range of plausible outputs. More recently, Wu et al. (2025) applied MBR to instructionfollowing tasks using an LLM-as-a-judge as a utility function. While this method yields strong results, it relies on a distillation step to approximate the utility, as directly querying an LLM judge during decoding is computationally prohibitive. In this work, we propose three lightweight adaptations to standard similarity-based utility functions, specifically designed for open-ended tasks characterised by high variability in latent structure.

Structure-Conditional Optimality 3

The central question addressed in this paper is how commonly employed utility functions for MBR decoding behave when complex structural variation is present. In Fig. 1, we illustrate the problem with a simplified example to highlight how the choice of utility function can influence decision-making-particularly when the outcome space contains multiple distinct modes. In this example, the outcome space is modelled as a bimodal Gaussian, and the decision problem is to select a single "best" outcome on the real line. If we use the negative squared error as our utility function,¹ the theoretical optimum corresponds to the mean of the bimodal distribution (the light blue line in Fig. 1). This solution may be undesirable as the mean lies in a region of low probability mass and is unlikely to be sampled in practice. If we apply a sampling-

¹Equivalently, one may frame this as minimising the risk under a squared error loss function.

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

based approximation to the decision rule, as is common in language generation applications of MBR, the approximation selects an outcome near this theoretical optimum, which typically resides at the boundary of one of the clusters. Alternatively, if we adopt a different utility function—such as a radial basis function kernel—the theoretical optimum shifts to the mode of the largest cluster (Fig. 1, dark blue line). This outcome, being more representative of a high-probability region, may be more desirable than either the low-probability intermodal mean or an outcome near the edge of a cluster.

258

259

260

262

263

264

267

269

271

272

273

275

276

277

278

281

284

290

292

296

297

301

304

In probability distributions over natural language, multiple such "modes" may also be present, albeit more difficult to define and detect. Generations might cluster around various semantically distinct plausible answers to a question, differing intended dialogue acts in a response, or varying response structure. Depending on the utility function used, this can result in behaviour analogous to that shown in Fig. 1. Whether this is desirable depends on the decision-maker; for instance, it may be appropriate if the model assigns probability mass to responses like The answer could be either [A] or [B], but in other cases, it could lead to suboptimal decisions. In this work, we investigate this phenomenon and propose simple adaptations to utility functions that encourage behaviour more similar to that of the RBF utility in the continuous example.

3.1 Evaluating Structural Sensitivity in MBR

To quantify the extent to which the MBR solution with commonly used utility functions respects structural variability in outcome spaces over natural language, we introduce two complementary metrics. These metrics evaluate whether MBR solutions align with, or differ from, solutions obtained when conditioning on latent structures.

Cluster Optimality. This metric quantifies the proportion of cases, over a test set, in which the MBR solution under the distribution P(Y|x)matches the MBR solution under the conditional distribution P(Y|x, s), where s denotes an annotated structure (e.g., a dialogue act) that we additionally condition on. Formally, let

$$\hat{y}_i = \operatorname*{argmax}_{h} \underset{P(Y|x_i)}{\mathbb{E}} [u(h, Y)]$$
(4)

be the MBR solution for input *i*, and

$$\hat{y}_i^{(s)} = \underset{h}{\operatorname{argmax}} \underset{P(Y|x_i,s)}{\mathbb{E}} [u(h,Y)]$$
(5)

the MBR solution conditioned on s. The cluster optimality metric is then defined, for test set D, as

$$\operatorname{CO} = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \mathbb{1}\{\hat{y}_i = \hat{y}_i^{(s)}\}$$
(6)

where $\mathbb{1}\{\cdot\}$ is the indicator function.

Cluster-Optimal Rank Correlation. In addition to the top-ranked solution, we also examine the full rankings produced by MBR. For each input *i*, consider a fixed set of hypothesis generations $\mathcal{H}_i = \{h_1, \ldots, h_n\}$ corresponding to structure *s*. Define the rankings:

$$R_{ij} = \text{rank of } h_j \text{ by } \underset{P(Y|x_i)}{\mathbb{E}} [u(h_j, Y)]$$
 (7)

$$R_{ij}^{(s)} = \text{rank of } h_j \text{ by } \underset{P(Y|x_i,s)}{\mathbb{E}} [u(h_j,Y)] \quad (8)$$

The cluster-optimal rank correlation is then the average Spearman's rank correlation coefficient ρ between these two rankings over the test set:

$$\operatorname{CORC} = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \rho(R_i, R_i^{(s)}) \tag{9}$$

4 Standard Utility Functions are Not Structure-Conditionally Optimal

We now demonstrate that MBR solutions derived using standard utility functions, such as BERTScore or BLEURT, often diverge from those obtained when conditioning on latent structural representations. While this divergence may be acceptable from the perspective of the decision-maker, our analysis adopts a different premise. We assume a production process in which the speaker first selects a latent structure-implicitly or explicitlyand then realises it through an utterance. Under this assumption, a generation should be optimal with respect to some latent structure, specifically the one selected by the speaker. Note that we do not model the initial stage of this process, *i.e.*, the selection or planning of the latent structure. Instead, we take it as given and focus on the requirement that the resulting generation be optimal within plausible realisations of the chosen structure.

To investigate how sensitive the MBR solution is to structural uncertainty in the outcome space, we consider three representative types of latent structure—dialogue act, emotion, and response structure—each of these defines a plausible axis of variation in generated text (see §4.1). For each structure, we construct a dataset that reflects the outcome space of a hypothetical model with uncertainty over its possible instantiations. We then compute the standard MBR solution over the entire space, and assess its optimality using the evaluation criteria introduced in §3.1. The results of this analysis are summarised in Tab. 1 and presented in §4.2.

349

353

354

357

361

363

364

381

4.1 Constructing Outcome Spaces with Controlled Structural Uncertainty

We ground our analysis in three types of latent structure. This section introduces each structure type and describes how we construct datasets to model uncertainty over their possible instantiations.

4.1.1 Latent Structure Types

We examine three types of latent structure that are representative of structural variability in the outcome spaces of open-ended generation tasks.

Dialogue Act. A dialogue act represents the communicative function or intent of an utterance within the context of a conversation. Following the taxonomy proposed by Amanova et al. (2016), we focus on four dialogue act types: INFORM, QUESTION, COMMISSIVE, and DIRECTIVE. Each of the four example utterances presented in the introduction exemplifies one of these categories.

Emotion. Another latent factor that shapes the
form of an utterance in conversation is the emotion
the speaker aims to express. In this work, we adopt
Ekman's six basic emotions (Ekman, 1992): HAPPINESS, SADNESS, FEAR, ANGER, SURPRISE, and
DISGUST. These emotional states influence both
lexical choice and broader stylistic features.

Response Structure. This structure type captures how information is organised within an instruction-following response. We consider four ad-hoc categories: BRIEF, a single-sentence reply; PARAGRAPH, a more developed, single-paragraph answer; LIST, a bullet-pointed set of items; and TABLE, a structured tabular presentation.

4.1.2 Dataset Construction

For each type of latent structure, we construct a dataset that simulates the outcome space of a hypothetical model with uncertainty over possible instantiations of that structure. As generation contexts, we randomly sample conversational contexts from the DailyDialog corpus (Li et al., 2017)— 1,000 each for dialogue act and emotion—and take the first 1,000 instructions from the Alpaca dataset

Metric	Utility	Dial. Act	Emotion	Resp. Str.	All (Avg)
СО	BERTScore BLEURT	0.370 0.410	0.330 0.510	0.390 0.530	0.363 0.483
CORC	BERTScore BLEURT	$\begin{array}{c} 0.081\\ 0.144\end{array}$	0.084 0.155	0.080 0.123	0.082 0.141

Table 1: Cluster Optimality (CO) and Cluster-Optimal Rank Correlation (CORC) of MBR solutions obtained using BERTScore and BLEURT utility functions over constructed outcome spaces.

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

(Taori et al., 2023) for response structure. We then prompt the instruction-tuned, 13 billion parameter variant² of the OLMo 2 model suite (OLMo et al., 2025) to generate outputs for each category within each structure type, using hand-curated prompts (see App. A for details). For every context, we generate 25 responses per structure category (*e.g.*, 25 BRIEF, 25 PARAGRAPH, 25 LIST, and 25 TA-BLE responses). This procedure results in 3,000 distinct outcome spaces, corresponding to 350,000 candidate generations in total.³

4.2 Structural Sensitivity of Standard MBR Utility Functions

Tab. 1 presents cluster optimality (CO, Eq. 6) and cluster-optimal rank correlation (CORC, Eq. 9) scores for MBR solutions under two standard utility functions across our three latent structure types. These metrics quantify how often the MBR-selected response is optimal with respect to the latent structure (CO), and how well it aligns with the structure-optimal ranking (CORC).

Across all structure types, we observe a consistent degree of suboptimality. The CO scores indicate that in fewer than half of the cases, the MBR solution is optimal with respect to its underlying structure (36.3% using BERTScore, 48.3% with BLEURT). This misalignment persists across dialogue act, emotion, and response structure, with no evident correlation to the number of clusters involved, suggesting that the failure to recover structure-optimal responses is not merely due to increased structural granularity. While slight differences are present between BLEURT and BERTScore, both utility functions consistently select suboptimal generations and demonstrate relatively weak ranking correlation. Overall, this analysis shows that standard utility functions possess low sensitivity to structural uncertainty.

²allenai/OLMo-2-1124-13B-Instruct

³The full dataset is included as Supplementary Material and will be released publicly upon acceptance.

480

481

482

5 Structure-Conditional MBR Decoding

To address the limitations of MBR decoding with 435 standard utility functions in the presence of struc-436 tural variability, we propose three structure-aware decoding approaches. 438

434

437

439

440

441

442

443

444

445

446

447

448

449

450

451

452

467

468

469

470

471

472

473

474

475

476

4

4

Utility Cut-off. Standard utility functions may implicitly penalise structural mismatches, but they do not prevent structurally dissimilar candidates from influencing the ranking of outputs. To mitigate this, we introduce a simple utility cut-off mechanism that filters out low-utility comparisons when computing expected utility. Specifically, we modify the utility function u(y, y') as follows:

$$u_{\rm cut}(y,y') = \begin{cases} u(y,y') & \text{if } u(y,y') \ge \tau, \\ \delta & \text{otherwise} \end{cases}$$
(10)

where τ is a threshold fixed across the dataset, and δ is a small constant (or zero). This limits the influence of distant or structurally irrelevant samples, aligning the MBR solution more closely with local modes in the outcome distribution.

Clustering. A more explicit approach to 453 structure-aware decoding is to first partition the 454 outcome space into clusters-each corresponding 455 to a distinct latent structure-and then apply MBR 456 within the dominant cluster. We implement this by 457 clustering candidate generations using sequence 458 embeddings $\phi(y)$ derived from a model fine-tuned 459 to detect particular structures of interest (e.g., dia-460 logue act, response structure, or affective content). 461 Formally, let $\mathcal{H} = \{h_1, \ldots, h_n\}$ be the set of 462 candidates, and let C_1, \ldots, C_k denote the resulting 463 clusters, with $\mathcal{H} = \bigcup_{j=1}^{k} C_j$. At inference time, 464 we restrict MBR decoding to the members of the 465 largest cluster $C^{\star} = \operatorname{argmax}_{C_i} |C_j|$ such that 466

$$\hat{y} = \operatorname*{argmax}_{h \in C^{\star}} \ \underset{P(Y|x)}{\mathbb{E}} [u(h,Y) \mid Y \in C^{\star}].$$
(11)

To recover a full ranking over candidates (e.g., for evaluation), we first rank clusters by size, and then rank candidates within each cluster based on expected utility. This two-stage process prioritises high-utility responses as judged against structurally consistent pseudo-references, reducing the risk of inter-modal averaging in the selected outputs.

> This procedure could also theoretically be formulated as an adaptation of the utility function:

77
$$u_{cl}(y, y') = \mathbb{1}\{C(y) = C(y')\} \times u(y, y') \times \mathbb{1}\{C(y) = C^*\}, \quad (12)$$

where C^* represents the cluster with highest probability mass under P(Y|x). Decoding then becomes standard MBR maximisation of expected utility under the adapted utility function.

Structure Embeddings. As an alternative to explicit clustering, we propose incorporating structural sensitivity into the utility function by leveraging structure-aware sequence embeddings. Specifically, we fine-tune a sequence embedding model to encode the structural property of interest and redefine the utility function to weight candidate comparisons by their similarity in this embedding space. Formally, for a candidate y and a reference y', we compute the modified utility as:

$$u_{\text{emb}}(y, y') = u(y, y') \cdot \cos\left(\phi(y), \phi(y')\right), \quad (13)$$

where u(y, y') is the original utility and $\cos(\cdot)$ denotes the cosine similarity between structuresensitive embeddings ϕ . To further reduce the influence of structurally mismatched samples, we also experiment with a threshold on cosine similarity: values below the threshold are set to zero, removing the contribution of the utility comparison to the expected utility altogether. This approach allows us to softly bias the MBR solution toward structurally coherent outputs, without requiring hard clustering or explicit structure labels at inference time.

6 **Experiments**

To evaluate the effectiveness of the proposed methods, we conduct a series of experiments on the dataset we constructed in §4, as well as two realworld instruction-following datasets. All our experiments use either BERTScore or BLEURT as the base utility function, two commonly employed utility functions in natural language generation (Freitag et al., 2022; Suzgun et al., 2023).

Cluster Optimality Under Controlled 6.1 **Structural Uncertainty**

We first assess our methods on the three datasets constructed in §4, which contain generations consisting of various types of structural uncertainty: over dialogue acts, emotions, and response structures. Recall that we treat these generations as hypothetical outcome spaces of a language model. That is, we consider all generations for a given context to be unbiased samples from a language model that we wish to perform MBR decoding with. We split the 1,000 contexts in each dataset into training, validation, and test sets using an 800/100/100 split.



Figure 2: Cluster Optimality and Cluster-Optimal Rank Correlation on the constructed outcome spaces of §4.1. We compare standard BERTScore and BLEURT MBR with the three adaptations to the utility functions proposed in §5.

Hyperparameter Selection. For each method proposed in §5, we use the training and validation splits to select hyperparameters and train the sequence embedding models. Determining the threshold in the Utility Cut-off approach is conducted separately for BERTScore and BLEURT, resulting in different thresholds. We base our sequence embedding models on the all-mpnet-base-v2⁴ Sentence Transformer (Reimers and Gurevych, 2019), which we further fine-tuned using a triplet loss and gold annotations of underlying structure to enhance sensitivity to the structural variation present in our datasets. We use the same sequence embedding models for our Clustering and Structure Embeddings approaches. We find that jointly fine-tuning and selecting thresholds on the combination of all three types of latent structure leads to the most robust performance in terms of CO,⁵ and we use these models for the experiments below. Further details on the hyperparameter selection and fine-tuning procedures can be found in App. B.

527

528

529

531

532

534

536

538

541

543

544

546

547

550

551

552

553

Results. We compare each of our proposed methods against standard sampling-based MBR decoding using either BERTScore or BLEURT as the utility function, and measure both cluster optimality (CO, Eq. 6) and cluster optimal ranking correlation (CORC, Eq. 9). Results are shown in

⁴https://huggingface.co/

sentence-transformers/all-mpnet-base-v2

Fig. 2. All the methods we proposed improve cluster optimality compared to the baseline versions of BERTScore and BLEURT. Utility Cutoff yields the smallest improvement over standard BERTScore and BLEURT MBR, on average increasing CO by 11.7% and 5.6%, respectively, and CORC by 0.091 and 0.002, respectively. The Clustering and Structure Embeddings approaches perform considerably better than baseline MBR. Clustering improves CO on average by 37.3%/27.7% and CORC by 0.382/0.320 CORC over standard BERTScore and BLEURT MBR, respectively. Similarly, Structure Embeddings improve CO on average by 38.7%/29.3% CO and CORC by 0.354/0.287 over standard BERTScore and BLEURT MBR, respectively. We note again that higher CO does not always correspond with higher CORC, indicating that achieving the clusteroptimal MBR solution is generally easier than recovering the entire ranking accurately. Additionally, we observe that some types of latent structure are more difficult to capture effectively than others.

554

555

556

557

558

559

561

562

563

564

565

566

567

568

569

570

571

573

574

575

576

577

578

579

580

581

582

584

6.2 Instruction-Following

Next, we evaluate our methods on two real-world instruction-following datasets: AlpacaEval (Li et al., 2023) and MT-Bench (Zheng et al., 2023). In this case, we do not have access to any labelling of potential latent structure. We use the same hyperparameters and sequence embedding models from the previous set of experiments, tuned on the combination of all three datasets from §4.

⁵Generally, we find CO and CORC in validation procedures to align reasonably well.

Benchmark	MBR	Cut-off	Cluster	Embeddings
AlpacaEval	96.5%	96.1%	97.0 %	96.1%
MT-Bench (single)	76.3%	90.0 %	80.0%	78.8%
MT-Bench (multi)	71.3%	70.0%	72.5%	74.4 %

Table 2: AlpacaEval and MT-Bench Prometheus win rates versus text-davinci-003 (AlpacaEval) / GPT-40 (MT-Bench). We compare standard BERTScore MBR with the approaches introduced in §5: Utility Cutoff, Clustering and Structure Embeddings.

As a language model, we select OLMo 2 (13B) (OLMo et al., 2025), and obtain 30 unbiased samples per prompt for use in MBR decoding. To measure task performance, we use Prometheus⁶ (Kim et al., 2024) as a judge, conducting relative grading against text-davinci-003 and GPT-40 (OpenAI, 2024) for AlpacaEval and MT-Bench, respectively.⁷ All experiments employ BERTScore as the base utility. Further details on the generation and evaluation procedures are provided in App. C.

586

587

588

591

592

Results. Tab. 2 reports win rates against text-davinci-003 and GPT-40 for standard 596 MBR decoding with a BERTScore utility, along-598 side our structure-conditional utilities from §5. On AlpacaEval, the Clustering method outperforms standard BERTScore MBR. In the single-turn MT-Bench setting, both Clustering and Utility Cut-off surpass standard MBR, with Utility Cut-off achiev-603 ing a notable 13.7 percentage point improvement and reaching a 90% win rate over GPT-40. This 604 indicates responses are often judged clearer, more helpful, accurate, and fully aligned with the intended purpose of the instruction. Performance declines across the board in the more challenging multi-turn MT-Bench setting. However, both clustering and structure embeddings still outperform 610 standard MBR, demonstrating improved structural 611 sensitivity also in extended interactions. Smaller 612 gains here may stem from reduced uncertainty 613 as conversational context accumulates, resulting 614 in less diverse outcome spaces. In such cases, 615 structure-conditional utilities likely yield results similar to standard MBR, reducing the relative ben-617 efit of structural adaptations. We also observe that 618 structure embeddings tend to outperform clustering, possibly because soft partitioning better captures 621 subtle structural differences, whereas hard clustering might inadvertently exclude partially similar 622

candidates. Nonetheless, the overall lower MBR performance in multi-turn tasks indicates these settings are inherently more challenging, beyond just reduced variability. Overall, the consistent improvements of structure-aware MBR methods over standard MBR suggest that incorporating latent structural information not only enhances the theoretical optimality of MBR solutions but also improves generation quality in practical settings. 623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

7 Conclusion

In this work, we examined the limitations of MBR decoding in open-ended generation scenarios, where outcome spaces might exhibit high structural variability. We hypothesised that commonly used utility functions are insufficiently sensitive to latent structural uncertainty, leading to suboptimal generation choices within structurally coherent clusters of responses. To evaluate this hypothesis, we constructed a dataset featuring naturally occurring contexts paired with outcome spaces that exhibit controlled variation in dialogue act, emotion, and response structure. Our findings confirm that MBR decoding under standard utilities frequently fails to select cluster-optimal candidates, with suboptimal selections occurring in more than half of the cases.

To address this issue, we proposed three approaches to adapt utility functions to be more structurally aware. The corresponding approaches— Clustering, Structure Embeddings, and Utility Cutoff—demonstrate significant improvements in both cluster optimality (up to 98% for response structure) and cluster-optimal rank correlation (up to 0.89 for response structure). Importantly, these methods incur only modest additional computational cost, requiring only lightweight fine-tuning of a sequence embedding model or performing a hyperparameter search for a threshold value.

Based on these results, we recommend adopting structure-aware MBR decoding in tasks characterised by medium to high outcome space variability, such as instruction-following and conversational tasks. We encourage future research into the development of structure-sensitive utility functions that build on this work to achieve even greater cluster optimality, generation quality, or inference-time efficiency. In addition, we see value in further investigating the relationship between outcome space variability and the effectiveness of structure-aware MBR, as well as the connection between cluster optimality and overall generation quality.

⁶prometheus-eval/prometheus-7b-v2.0

⁷We did not find any available multi-turn system generations for the full MT-Bench dataset. Therefore, we generated our own from OpenAI's GPT-40 using greedy decoding.

695

701

703

704

705

706

710

711

712

713

714

715

716

717

718

719

721

Limitations

To test our hypothesis on the sub-optimality of 674 standard similarity-based MBR utility functions, 675 we relied on a curated dataset that captures three representative types of latent structure commonly found in open-ended natural language generation tasks. However, this dataset does not exhaustively 679 cover all possible structural variations present in natural language. Additionally, our evaluation assumes that language models accurately represent uncertainty over latent structures-an assumption that may not always hold in practice (see, e.g., Giulianelli et al., 2023). For example, in a dialogue setting, a model might assign most of its probability mass to responses aligned with the INFORM category, even if human responses display a broader range of structural types.

> In terms of computational requirements, our methods introduce minimal overhead beyond standard MBR decoding. However, it is important to note that MBR decoding itself *is* significantly more computationally demanding than greedy decoding or sampling a single generation. Since our approaches build on MBR, they inherit this higher computational cost. That said, we believe our methods stand to benefit from recent advances aimed at improving the efficiency of MBR decoding (Cheng and Vlachos, 2023; Vamvas and Sennrich, 2024; Yang et al., 2024).

Finally, in our evaluation on instructionfollowing datasets, we rely on an LLM-as-a-judge: Prometheus (Kim et al., 2024). These are imperfect evaluators, may be biased towards particular types of responses (Wang et al., 2024; Stureborg et al., 2024), such as more elaborate ones, and do not always align with human judgements (Zeng et al., 2024; Bavaresco et al., 2024). Additionally, Prometheus relies on a predefined rubric, and its performance may be sensitive to the specific formulation of that rubric. We did not conduct extensive experiments with alternative rubric designs, which may influence the robustness of the results.

References

Dilafruz Amanova, Volha Petukhova, and Dietrich Klakow. 2016. Creating annotated dialogue resources: Cross-domain dialogue act classification. In Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16), pages 111–117, Portorož, Slovenia. European Language Resources Association (ELRA).

Joris Baan, Nico Daheim, Evgenia Ilia, Dennis Ulmer, Haau-Sing Li, Raquel Fernández, Barbara Plank, Rico Sennrich, Chrysoula Zerva, and Wilker Aziz. 2023. Uncertainty in natural language generation: From theory to applications. 723

724

725

726

727

728

729

730

731

732

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

778

- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, André F. T. Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2024. LLMs instead of human judges? A large scale empirical study across 20 NLP evaluation tasks.
- James O Berger. 1985. *Statistical decision theory and Bayesian analysis; 2nd ed.* Springer Series in Statistics. Springer, New York.
- Julius Cheng and Andreas Vlachos. 2023. Faster minimum Bayes risk decoding with confidence-based pruning. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12473–12480, Singapore. Association for Computational Linguistics.
- Julius Cheng and Andreas Vlachos. 2024. Measuring uncertainty in neural machine translation with similarity-sensitive entropy. In Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2115–2128, St. Julian's, Malta. Association for Computational Linguistics.
- Yuntian Deng, Volodymyr Kuleshov, and Alexander Rush. 2022. Model criticism for long-form text generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11887–11912, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2020. Is MAP decoding all you need? the inadequacy of the mode in neural machine translation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2022. Sampling-based approximations to minimum Bayes risk decoding for neural machine translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- P Ekman. 1992. Are there basic emotions? *Psychological review*, 99(3):550–553.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of*

890

891

836

837

the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1396–1412, Seattle, United States. Association for Computational Linguistics

779

783

789

790

791

793

794

797

798

800

801

803

804

809

810

811

812

813

815

817

818

819

821

825

827

829

830

831

834

- Seattle, United States. Association for Computational Linguistics.Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model
- Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.
- Mario Giulianelli, Joris Baan, Wilker Aziz, Raquel Fernández, and Barbara Plank. 2023. What comes next? evaluating uncertainty in neural text generators against human production variability. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14349–14371, Singapore. Association for Computational Linguistics.
 - Evgenia Ilia and Wilker Aziz. 2024. Variability need not imply error: The case of adequate but semantically distinct responses.
- Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. 2024. Prometheus 2: An open source language model specialized in evaluating other language models.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. 2024.
 Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation.
 In *The Eleventh International Conference on Learning Representations*.
- Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004, pages 169–176, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpacaeval: An automatic evaluator of instructionfollowing models. https://github.com/ tatsu-lab/alpaca_eval.
- Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. In *Proceedings* of The 8th International Joint Conference on Natural Language Processing (IJCNLP 2017).
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, Nathan Lambert, Dustin Schwenk, Oyvind Tafjord, Taira Anderson, David Atkinson, Faeze Brahman, Christopher Clark, Pradeep Dasigi, Nouha Dziri, Michal

Guerquin, Hamish Ivison, Pang Wei Koh, Jiacheng Liu, Saumya Malik, William Merrill, Lester James V. Miranda, Jacob Morrison, Tyler Murray, Crystal Nam, Valentina Pyatkin, Aman Rangapur, Michael Schmitz, Sam Skjonsberg, David Wadden, Christopher Wilhelm, Michael Wilson, Luke Zettlemoyer, Ali Farhadi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. 2 olmo 2 furious.

- OpenAI. 2024. Gpt-4o system card.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702, Online. Association for Computational Linguistics.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- Felix Stahlberg and Bill Byrne. 2019. On NMT search errors and model errors: Cat got your tongue? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3356– 3362, Hong Kong, China. Association for Computational Linguistics.
- Rickard Stureborg, Dimitris Alikaniotis, and Yoshi Suhara. 2024. Large language models are inconsistent and biased evaluators.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. 2023. Follow the wisdom of the crowd: Effective text generation via minimum Bayes risk decoding. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4265–4293, Toronto, Canada. Association for Computational Linguistics.
- Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Stanford alpaca: An instruction-following llama model. https://github.com/tatsu-lab/ stanford_alpaca.
- Jannis Vamvas and Rico Sennrich. 2024. Linear-time minimum Bayes risk decoding with reference aggregation. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics*

- _

5-

(Volume 2: Short Papers), pages 790–801, Bangkok, Thailand. Association for Computational Linguistics.

- Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and Zhifang Sui. 2024. Large language models are not fair evaluators. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 9440–9450, Bangkok, Thailand. Association for Computational Linguistics.
- Ian Wu, Patrick Fernandes, Amanda Bertsch, Seungone Kim, Sina Khoshfetrat Pakazad, and Graham Neubig.
 2025. Better instruction-following through minimum bayes risk. In *The Thirteenth International Conference on Learning Representations.*
- Guangyu Yang, Jinghong Chen, Weizhe Lin, and Bill Byrne. 2024. Direct preference optimization for neural machine translation with minimum Bayes risk decoding. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers), pages 391–398, Mexico City, Mexico. Association for Computational Linguistics.
- Zhiyuan Zeng, Jiatong Yu, Tianyu Gao, Yu Meng, Tanya Goyal, and Danqi Chen. 2024. Evaluating large language models at evaluating instruction following. In *The Twelfth International Conference on Learning Representations*.
- Hugh Zhang, Daniel Duckworth, Daphne Ippolito, and Arvind Neelakantan. 2021. Trading off diversity and quality in natural language generation. In *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 25–33, Online. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging Ilm-as-a-judge with mt-bench and chatbot arena.

A Data Generation

In §4.1.2 we semi-automatically construct a natural language dataset of hypothetical outcome spaces with varying underlying latent structure: dialogue act, emotion and response structure. Here, we describe the full generation procedure in more detail.

A.1 Generation Procedure

We use the input prompts from the DailyDialog dataset (Li et al., 2017) as a basis for the dialogue

act and emotion subsets, and Alpaca (Taori et al., 2023) as a basis for the response structure subset. In order to obtain sampled generations, we use allenai/OLMo-2-1124-13B-Instruct (OLMo et al., 2025).

- Pre-processing. For DailyDialog, we remove unnecessary white spaces, and filter dialogues by the number of turns (only dialogues with ≥ 2 turns are kept). For each of the dialogues, we also randomly sample the number of turns which we will use from the dialogue, between two and the total number of available turns minus two. We also prepend "A:" or "B:" prefixes to each turn to indicate the speaker. For Alpaca, we only filter out any prompts that have an additional input field.
- 2. **Prompt creation**. We randomly select 1,000 dialogues twice from the DailyDialog dataset and the first 1,000 instructions without input field from the Alpaca dataset. We then fill out the predefined prompt templates (defined in App. A.2) with the selected examples, resulting in prompts for each pre-defined category within every latent structure. In total, we create 1,000 * 4 = 4,000 inputs for the dialogue act subset, 1,000 * 4 = 4,000 inputs for response structure subset, and 1,000 * 6 = 6,000 inputs for the emotion subset to be fed into OLMo.
- 3. **Generation**. Using OLMo, we then generate 25 unbiased samples for each of the 14,000 constructed input prompts.
- 4. **Post-processing**. We finalize the procedure by removing quotation marks around the generation and stripping any "*A*:" or "*B*:" prefixes at the start of generations.

A.2 Prompt Templates

To generate prompts covering all categories within a latent structure, we define three types of prompt templates - one per latent structure.

A.2.1 Dialogue Act

For the dialogue act subset, we define a responder as the speaker whose turn it is now to speak (*e.g.* if the dialogue excerpt ends at A's turn, we define responder = B). For each dialogue from DailyDialog we then iterate through the four defined dialogue acts and pass each of them as act_name.

993	
994	### **Types of Dialoque Acts**
995	Here are common categories of dialogue
996	acts, though exact categorizations
997	may vary depending on the framework:
998	#### **1 Inform**
999	- The Inform class contains all
1000	statements and questions by which
1001	the speaker shares information
1001	with the listoper. The speaker
1002	with the instenet. The speaker
1003	assumes the information is
1004	correct and believes the
1005	addressee does not know or is not
1006	aware of it yet.
1007	- **Examples**:
1008	- "The meeting starts at 3 PM."
1009	- "I've already emailed the report."
1010	- "I saw John at the store
1011	yesterday."
1012	– "William Shakespeare wrote it."
1013	- "The Eiffel Tower is in Paris."
1014	
1015	#### **2. Question**
1016	- The Question class includes speech
1017	acts where the speaker seeks
1018	information by asking a question.
1019	These acts are used when the
1020	speaker wants to know something
1021	and believes the listener has the
1022	answer. Questions can take
1023	different forms, including
1024	Propositional Questions (yes/no
1025	questions), Check Questions
1026	(confirming known information),
1027	Set Questions (open-ended
1028	questions), and Choice Questions
1029	(questions with multiple options)
1030	- **Examples***
1031	- "Did you finish your assignment?"
1032	- "You've met Sarah before haven't
1033	vou?"
1034	- "What time does the meeting
1035	start?"
1035	- "Could you clarify what you meant
1030	by that?"
1037	- "Do you profor coffee or too?"
1020	Do you pierer corree or cea:
1039	#### 2 Directive
1040	#### **3. Directive **
1041	- The Directive class includes speech
1042	listenen to perform an action
1043	This sleep course Demosts
1044	(achier corrects de corrects)
1045	(asking someone to do something),
1046	Instructions (giving direct
1047	orders or guidance), Suggestions
1048	(offering recommendations), and
1049	Accepting or Rejecting Offers
1050	(responding to proposals). These
1051	acts differ based on how much
1052	pressure the speaker applies and
1053	their assumptions about the
1054	listener's willingness and
1055	ability to comply.
1056	- **Examples**:
1057	- "Can you send me the file?"
1058	- "Fill out this form before the
1059	appointment."
1060	- "You should try the new Italian
1061	restaurant downtown."

- "Yes, I'd love to join you for	1062
dinner!"	1063
 "No, I can't take on another 	1064
project right now."	1065
	1066
#### **4. Commissive**	1067
- The Commissive class involves	1068
speech acts where the speaker	1069
commits to performing an action	1070
in the future. These acts include	1071
Accepting or Rejecting Requests,	1072
Suggestions, and Offers. By	1073
performing a Commissive act, the	1074
speaker is making a promise or	1075
commitment to carry out the	1076
action requested, suggested, or	1077
offered. These acts reflect the	1078
speaker's willingness to take	1079
responsibility for fulfilling the	1080
commitment, whether by agreeing	1081
to a proposal or refusing it.	1082
- **Examples**:	1083
- "Fine, I'll pick you up at 5 PM."	1084
 "Sorry, I can't do that right 	1085
now."	1086
 "That sounds great, I'll take the 	1087
promotion."	1088
- "I promise to finish the report	1089
by the end of the day."	1090
 "I'll make sure to take care of 	1091
it this weekend."	1092
	1093
	1094
	1095
### **Dialogue Excerpt**	1096
	1097
{dialogue from DailyDialog}	1098
	1099
	1100
	1101
### **Instructions**	1102
Please consider the provided dialogue	1103
excerpt and provide a plausible	1104
response (and only a single	1105
response) for {responder} that	1106
reflects the following dialogue	1107
<pre>act: {act_name}. Output only</pre>	1108
{responder}'s response with no	1109
additional text. <end_of_prompt></end_of_prompt>	1110
A.2.2 Emotion	1111
For the emotion subset, we again define a	1112

For the emotion subset, we again define a responder in the same way as in App. A.2.1. For each dialogue from DailyDialog we then iterate through the six defined emotions and pass each of them as emotion_name.

```
### **Types of Emotions**
                                                 1117
Here are seven main categories of
                                                 1118
   emotions.
                                                 1119
                                                 1120
#### **1. Anger**
                                                 1121
  - The Anger category represents
                                                 1122
                                                 1123
     emotions related to feelings of
     displeasure, hostility, or
                                                 1124
     frustration. This emotion often
                                                 1125
     arises when someone feels wronged
                                                 1126
                                                 1127
      or blocked from achieving their
```

1113

1114

1115

1116

1128	goal. It can range from mild
1120	irritation to intense rade
1120	- ++Fyamples++.
1131	- "I can't believe this is
1122	happoning.
1132	"This is so unfairl"
113/	- "Why doos overything always go
1135	- Why does everything always go
1100	WIDING IDI IIIE:
1127	- "I'M SO IIUSLIALEU WILH LHIS
1122	- "I'm really mad about how things
1120	- "I'll really had about now things
11/0	curned out."
1140	
1141	#### **2. DISGUSL**
1142	- The Disgust Category Includes
1143	emotions related to a strong
1144	sense of revulsion, disapproval,
1145	or distaste. It often arises when
1146	something is perceived as
1147	offensive, repellent, or morally
1148	objectionable.
1149	- **Examples**:
1150	- "That food looks awful!"
1151	- "I can't stand how they treat
1152	people."
1153	- "This is disgusting. I can't
1154	believe they did that."
1155	- "I feel sick just thinking about
1156	it."
1157	- "That's absolutely revolting!"
1158	
1159	#### **3. Fear**
1160	- The Fear category includes emotions
1161	related to anxiety, nervousness,
1162	and concern about possible danger
1163	or harm. Fear can be rational or
1164	irrational and may cause physical
1165	or psychological distress.
1166	- **Examples**:
1167	 "I'm really scared about what's
1168	going to happen."
1169	- "I don't know if I can handle
1170	this situation."
1171	– "What if things don't go as
1172	planned?"
1173	- "I'm afraid something bad might
1174	happen."
1175	- "I'm nervous about the meeting
1176	this morning."
1177	
1178	#### **4. Happiness**
1179	- The Happiness category includes
1180	emotions related to joy,
1181	contentment, and pleasure.
1182	Happiness is often associated
1183	with positive experiences,
1184	accomplishments, and satisfying
1185	events.
1186	- **Examples**:
1187	- "I'm so excited about this
1188	weekend!"
1189	- "This is such a great day!"
1190	- "I feel so happy about my
1191	progress."
1192	- "That sounds amazing, I'm really
1193	looking forward to it!"
1194	 "I'm so glad everything worked
1195	out!"
1196	
1197	#### **5. Sadness**

- The Sadness category represents	1198
emotions related to feelings of	1199
loss, disappointment, or sorrow.	1200
It often arises when there is a	1201
sense of unmet expectations,	1202
failure, or grief.	1203
- **Examples**:	1204
- "I feel so down about what	1205
happened."	1206
- "I can't stop thinking about it,	1207
it's just so upsetting."	1208
- "I'm really sad things turned out	1209
this way."	1210
- "It's been a tough time, and I	1211
feel heartbroken."	1212
 "I don't know how to get over 	1213
this sadness."	1214
	1215
#### **6. Surprise**	1216
- The Surprise category represents	1217
emotions related to unexpected	1218
events or outcomes, ranging from	1219
shock to awe. This emotion can be	1220
positive or negative, depending	1221
on the nature of the surprise.	1222
- **Examples**:	1223
- "Wow, I didn't see that coming!"	1224
- "That's such a surprise, I can't	1225
believe it!"	1226
 "I'm totally shocked by what 	1227
happened."	1228
- "I wasn't expecting that at all!"	1229
- "I'm so surprised you did that!"	1230
	1231
	1232
	1233
### **Dialogue Excerpt**	1234
	123
{dialogue from DailyDialog}	1236
	1237
	1238
	1239
### **Instructions**	1240
Please consider the provided dialogue	1241
excerpt and provide a plausible	1242
response (and only a single	1243
response) for {responder} that	1244
reflects the following emotion:	124
{emotion_name}. Output only	1240
{responder}'s response with no	124/
additional text. <end_oi_prompt></end_oi_prompt>	1240

A.2.3 Response Structure

For the response structure subset, we define four different prompt templates, one per category of response structure. For each prompt from Alpaca we then append each of these templates, resulting in four different prompts, one per category, per input instruction.

BRIEF

```
{prompt from Alpaca} Give me a brief 1258
  sentence with the answer. Make 1259
  sure to restrict your response 1260
  to a single sentence. 1261
```

1249

1250

1251

1252

1253

1254

1255 1256

1257

```
1262
1263
1264
1265
1266
```

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1281

1282

1283

1284

1285

1286

1287

1288

1289

1291

1292

1293

1294

1295

1296

1297

1298

1299

PARAGRAPH

{prompt from Alpaca} W	rite an
extensive paragraph	n on the
topic. Restrict you	ir answer to a
single paragraph	

List

{prompt from Alpaca} In your answer, make sure to include a bullet point list of items relevant to the topic. Keep your answer brief and make sure it contains a bullet point list.

TABLE

{prompt from Alpaca} In your answer, include a table relevant to the topic. Keep your answer brief and make sure it contains a table.

B Hyperparameter Selection

We randomly split our generated datasets (dialogue act, emotion and response structure) into 800/100/100 training/validation/testing data points. All data consists of multiple labelled clusters with 25 generations per cluster. We compute BERTScore and BLEURT MBR solutions conditioned on each labelled cluster to get clusteroptimal rankings and MBR solutions to compare to. We use the training and validation splits for fine-tuning sequence embedding models and for hyperparameter selection. We have performed all training and hyperparameter selection both on individual datasets (either dialogue act, emotion or response structure) as well as on the combination of all datasets. We find that models trained on all data perform best overall and thus have used those in experiments. We proceed here to discuss the results of hyperparameter selection for each individual approach in more detail.

1300 **Utility Cut-off.** We considered both an absolute threshold on the utility value as well as a threshold 1301 on the deviation from the highest observed utility 1302 in the sample. We don't consider any utility com-1303 parisons with the candidate itself, *i.e.*, we mask out 1304 the diagonal of the utility matrix. Furthermore, we 1305 experiment with both setting utility values below 1306 the threshold to 0 or -1 as well as discarding those utility comparisons altogether. We test a range of 1308 50 threshold values ranging within reasonable val-1309 ues for the utility function itself, and order settings 1310 based on cluster optimality on the training data. We 1311 then take the 10 best performing setups and select 1312

the one that performs best in terms of cluster opti-1313 mality on the validation data. We tune the threshold 1314 independently for both BLEURT and BERTScore. 1315 We find an absolute value threshold to work supe-1316 rior for both utilities, as well as to zero out values 1317 below the threshold. We find an optimal threshold 1318 of 0.512 and 0.918 for BLEURT and BERTScore 1319 respectively. 1320

1321

1322

1323

1324

1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

1360

1361

1362

We use the Sentence Transformers Clustering. all-mpnet-base-v2 model as a basis for obtaining sequence embeddings. We further fine-tune this model using a triplet loss on triplets from our labelled datasets. We experiment with learning rates between 1×10^{-4} , 1×10^{-5} and 1×10^{-6} and find a learning rate of 1×10^{-5} to lead to best validation loss overall. We then use these sequence embeddings with the k-means algorithm to obtain clusters. We select a number of clusters based on the silhouette score for k = [2, 6] and set a threshold that the silhouette scores need to reach, otherwise k is set to 1 and we consider all generations to come from a single cluster. This threshold is tuned on prediction accuracy on the number of clusters for a range of values in (0, 1) on random subsamples of the validation data containing a random number of clusters per subsample.

Structure Embeddings. Here, we use the same fine-tuned sentence transformer model from the clustering approach. We shift and compress cosine similarity values to range between 0 and 1. We optionally consider a threshold on cosine similarity and perform an identical selection procedure to that for the threshold in the Utility Cut-off approach. We find that a threshold does considerably improve cluster optimality and end up with a threshold of 0.918.

Fine-Tuned Utilities: BERTScore and BLEURT. We also attempted fine-tuning BERTScore and BLEURT directly to be more sensitive to the latent structures we expect in the data. We experimented with fine-tuning BERTScore with a triplet loss on the sequence embeddings of the underlying roberta-large model, and used a mean squared error regression loss to fine-tune BLEURT to predict comparisons with out-of-cluster generations as 0 or -1. We attempted a range of hyperparameter values, but found that the resulting utility functions performed poorly across the board. Hence, we have not included those models in the main paper. You are an advanced AI assistant specializing in clear, well-reasoned, and articulate responses. Your goal is to provide comprehensive and accurate answers while ensuring coherence, logical consistency, and factual correctness. Be precise, provide evidence-based explanations, and use structured reasoning when appropriate. If a question has multiple interpretations, clarify them before answering. Avoid unnecessary verbosity while maintaining completeness. If uncertain, state your level of confidence and explain why.

Figure 3: System prompt provided to the model for all decoding methods when generating for AlpacaEval and MT-Bench.

1363 1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1383

1384

1385

1386

1387

1388

1389

C Evaluation on AlpacaEval and MT-Bench

We performed evaluation on AlpacaEval and MT-Bench using Prometheus as an LLM-as-a-judge model.⁸ We test the instruction-following capabilities as follows:

- 1. **Generation**. For each instruction from the dataset, we generate our answers from each respective decoding method. We use a system prompt given in Fig. 3. When generating for single-turn MT-Bench, we only prompt the model with the first turn and save its output. When generating for multi-turn MT-Bench, we first prompt the model with just the first turn, save its output, and then we prompt it again with both turns and the reference GPT-40 generation to the first prompt.⁹ Therefore, the total number of instructions for multi-turn MT-Bench.
 - Evaluation. We then pass the instruction, the reference answer, as well as the generations of our decoding methods to Prometheus. The reference answers for AlpacaEval (included in the dataset) were generated by text-davinci-003. We collected reference answers for multi-turn MT-Bench our-

selves doing greedy decoding from the GPT-1390 40 through the OpenAI API. We use the 1391 pre-defined RELATIVE_PROMPT_WO_REF 1392 prompt template for Prometheus, as well as 1393 relative grading - for each pair of competing 1394 outputs, Prometheus returns one letter (A or 1395 B) defining which output is preferred. We de-1396 fine the rubric as "Is the answer clear, helpful, 1397 accurate, and fully aligned with the intended 1398 purpose of the instruction?" 1399

Final Score. For every decoding method we then calculate that methods' win rate over the set of reference generations according to Prometheus. In multi-turn MT-Bench, we report the average of the win rates of both turns.

 $^{^{8}}We$ use prometheus-eval/prometheus-7b-v2.0 ^{9}We opted to always provide the reference response in multi-turn MT-Bench, as to not suffer from compounding errors.