

## A Appendix: Inappropriate Gaussian Priors

The majority of work on VAEs propose to parametrize  $z$  — both the prior and approximate posterior (encoder) — as a multivariate Gaussian variable. However, the multivariate Gaussian is a uni-modal distribution and can therefore only represent one mode in latent space. Furthermore, the multivariate Gaussian is perfectly symmetric with a constant kurtosis. These properties are problematic if the latent variables we aim to represent are inherently multi-modal, or if the latent variables follow complex, non-linear probability manifolds (e.g. asymmetric distributions or heavy-tailed distributions). For example, the frequency of topics in news articles could be represented by a continuous probability distribution, where each topic has its own island of probability mass; *sports* and *politics* topics might each be clustered on their own separate island of probability mass with zero or little mass in between them. Due to its uni-modal nature, the Gaussian distribution can never represent such probability distributions. As another example, ambiguity and uncertainty in natural language conversations could similarly be represented by islands of probability mass; given the question *How do I install Ubuntu on my laptop?*, a model might assign positive probability mass to specific, unambiguous entities like *Ubuntu 4.10* and to well-defined procedures like *installation using a DVD*. In particular, certain entities like *Ubuntu 4.10* are now outdated — these entities occur rarely in practice and should be considered rare events. When modeling such complex, multi-modal latent distributions, the mapping from multivariate Gaussian latent variables to outputs — i.e. the conditional distribution  $P_\theta(w_n|z)$  — has to be highly non-linear in order to compensate for the simplistic Gaussian distribution and capture the natural latent factors in an intermediate layer of the model. However, it is difficult to learn such non-linear mappings when using the variational bound in eq. (1), as it incurs additional variance from sampling the latent variable  $z$ . Consequently, such models are likely to converge on solutions that do not capture salient aspects of the latent variables, which in turn leads to a poor fit of the output distribution.

## B Appendix: Piecewise Constant Variable Derivations

To train the model using the re-parametrization trick, we need to generate  $z = f(\epsilon)$  where  $\epsilon \sim \text{Uniform}(0, 1)$ . To do so, we employ inverse transform sampling (Devroye, 1986), which requires finding the inverse of the cumulative distribution function (CDF). We derive the CDF of eq. (2):

$$\phi(z) = \frac{1}{K} \sum_{i=1}^n \mathbb{1}\left(\frac{i}{n} \leq z\right) K_i + 1 \left(\frac{i-1}{n} \leq z \leq \frac{i}{n}\right) * \left(z - \frac{i-1}{n}\right) a_i. \quad (3)$$

Next, we derive its inverse:

$$\phi^{-1}(\epsilon) = \sum_{i=1}^n \mathbb{1}\left(\frac{1}{K} \sum_{j=0}^{i-1} K_j \leq \epsilon \leq \frac{1}{K} \sum_{j=0}^i K_j\right) * \left(\frac{i-1}{n} + \frac{K}{a_i} \left(\epsilon - \frac{1}{K} \sum_{j=0}^{i-1} K_j\right)\right) \quad (4)$$

Armed with the inverse CDF, we can now draw a sample  $z$ :

$$z = \phi^{-1}(\epsilon), \quad \text{where } \epsilon \sim \text{Uniform}(0, 1). \quad (5)$$

In addition to sampling, we need to compute the Kullback-Leibler (KL) divergence between the prior and approximate posterior distributions of the piecewise constant variables. We assume both the prior and the posterior are piecewise constant distributions. We use the *prior* superscript to denote prior parameters and the *post* superscript to denote posterior parameters (encoder model parameters). The KL divergence between the prior and posterior can be computed using a sum of integrals, where each integral inside the sum corresponds to one constant segment:

$$\begin{aligned} & \text{KL}[Q_\psi(z|w_1, \dots, w_N) || P_\theta(z)] \\ &= \int_0^1 Q_\psi(z|w_1, \dots, w_N) \log\left(\frac{Q_\psi(z|w_1, \dots, w_N)}{P_\theta(z)}\right) dz \end{aligned} \quad (6)$$

$$= \sum_{i=1}^n \int_0^{1/n} \frac{a_i^{\text{post}}}{K^{\text{post}}} \log\left(\frac{a_i^{\text{post}}/K^{\text{post}}}{a_i^{\text{prior}}/K^{\text{prior}}}\right) dz \quad (7)$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{a_i^{\text{post}}}{K^{\text{post}}} \log\left(\frac{a_i^{\text{post}}/K^{\text{post}}}{a_i^{\text{prior}}/K^{\text{prior}}}\right) \quad (8)$$

$$\begin{aligned} &= \frac{1}{n} \frac{1}{K^{\text{post}}} \sum_{i=1}^n a_i^{\text{post}} \left(\log(a_i^{\text{post}}) - \log(a_i^{\text{prior}})\right) \\ &+ \log(K^{\text{prior}}) - \log(K^{\text{post}}) \end{aligned} \quad (9)$$

In order to improve training, we further transform the piecewise constant latent variables to lie within the interval  $[-1, 1]$  after sampling:  $z' = 2z - 1$ . This ensures the input to the decoder RNN has mean zero initially.

### C Appendix: NVDM Implementation

The complete NVDM architecture is defined as:

$$\begin{aligned}\pi(W) &= f^0(E^0W + b^0), \\ \text{Enc}(W) &= f^1(E^1\pi(W) + b^1), \\ z_{\text{Gaussian}} &= \mu^{\text{post}} + \sqrt{\sigma^{2,\text{post}}} \otimes \epsilon_0, \\ z_{\text{Piecewise}} &= \phi^{-1,\text{post}}(\epsilon_1), \\ z &= \langle z_{\text{Gaussian}}, z_{\text{Piecewise}} \rangle, \\ \text{Dec}(w, z) &= g(-w^T Rz),\end{aligned}$$

where  $\otimes$  is the Hadamard product,  $\langle \circ, \circ \rangle$  is an operator that combines the Gaussian and the Piecewise variables and  $\text{Dec}(w, z)$  is the decoder model.<sup>5</sup> As a result of using the re-parametrization trick and choice of prior, we calculate the latent variable  $z$  through the two samples,  $\epsilon_0$  and  $\epsilon_1$ .  $f(\circ)$  is a non-linear activation function, which was the parametrized linear rectifier (with a learnable ‘‘leak’’ parameters) for the 20 News-Groups experiments and the softsign function, or  $f(v) = v/(1 + |v|)$ , for Reuters and CADE. The decoder model  $\text{Dec}(z)$  outputs a probability distribution over words conditioned on  $z$ . In this case, we define  $g(\circ)$  as the softmax function (omitting the bias term  $c$  for clarity) computed as:

$$\text{Dec}(w, z) = P_\theta(w|z) = \frac{\exp(-w^T Rz)}{\sum_{w'} \exp(-w'^T Rz)},$$

The decoder’s output is used to calculate the first term in the variational lower-bound:  $\log P_\theta(W|z)$ . The prior and posterior distributions are used to compute the KL term in the variational lower-bound. The lower-bound is:

$$\begin{aligned}\mathcal{L} &= \mathbb{E}_{Q_\psi(z|W)} \left[ \sum_{i=1}^N \log P_\theta(w_i|z) \right] \\ &\quad - \text{KL}[Q_\psi(z|W) || P_\theta(z)],\end{aligned}$$

where the KL term is the sum of the Gaussian and piecewise KL-divergence measures:

$$\begin{aligned}\text{KL}[Q(z|W) || P(z)] &= \text{KL}_{\text{Gaussian}}[Q(z|W) || P(z)] \\ &\quad + \text{KL}_{\text{Piecewise}}[Q(z|W) || P(z)].\end{aligned}$$

<sup>5</sup>Operations include vector concatenation, summation, or averaging.

The KL-terms may be interpreted as regularizers of the parameter updates for the encoder model (Kingma and Welling, 2014). These terms encourage the posterior distributions to be similar to their corresponding prior distributions, by limiting the amount of information the encoder model transmits regarding the output.

### D Appendix: VHRED Implementation

As described in the model section, the probability distribution of the generative model factorizes as:

$$\begin{aligned}P_\theta(\mathbf{w}_1, \dots, \mathbf{w}_N) &= \prod_{n=1}^N P_\theta(\mathbf{w}_n | \mathbf{w}_{<n}, z_n) P_\theta(z_n | \mathbf{w}_{<n}) \\ &= \prod_{n=1}^N \prod_{m=1}^{M_n} P_\theta(w_{n,m} | w_{n,<m}, \mathbf{w}_{<n}, z_n) P_\theta(z_n | \mathbf{w}_{<n}),\end{aligned}\tag{10}$$

where  $\theta$  are the model parameters. VHRED uses three RNN modules: an *encoder* RNN, a *context* RNN and a *decoder* RNN. First, each utterance is encoded into a vector by the *encoder* RNN:

$$\begin{aligned}h_{n,0}^{\text{enc}} &= \mathbf{0}, \quad h_{n,m}^{\text{enc}} = f_\theta^{\text{enc}}(h_{n,m-1}^{\text{enc}}, w_{n,m}) \\ &\quad \forall m = 1, \dots, M_n,\end{aligned}$$

where  $f_\theta^{\text{enc}}$  is either a GRU or a bidirectional GRU function. The last hidden state of the *encoder* RNN is given as input to the *context* RNN. The *context* RNN uses this state to update its internal hidden state:

$$h_0^{\text{con}} = \mathbf{0}, \quad h_n^{\text{con}} = f_\theta^{\text{con}}(h_{n-1}^{\text{con}}, h_{n,M_n}^{\text{enc}}),$$

where  $f_\theta^{\text{con}}$  is a GRU function taking as input two vectors. This state conditions the prior distribution over  $z_n$ :

$$P_\theta(z_n | \mathbf{w}_{<n}) = f_\theta^{\text{prior}}(z_n; h_{n-1}^{\text{con}}),\tag{11}$$

where  $f^{\text{prior}}$  is a PDF parametrized by both  $\theta$  and  $h_{n-1}^{\text{con}}$ . Next, a sample is drawn from this distribution:  $z_n \sim P_\theta(z_n | \mathbf{w}_{<n})$ . The sample and *context* state are given as input to the *decoder* RNN:

$$\begin{aligned}h_{n,0}^{\text{dec}} &= \mathbf{0}, \quad h_{n,m}^{\text{dec}} = f_\theta^{\text{dec}}(h_{n,m-1}^{\text{dec}}, h_{n-1}^{\text{con}}, z_n, w_{n,m}) \\ &\quad \forall m = 1, \dots, M_n,\end{aligned}$$

where  $f_\theta^{\text{dec}}$  is the LSTM gating function taking as input four vectors. The output distribution is computed by passing  $h_{n,m}^{\text{dec}}$  through an MLP  $f_\theta^{\text{mlp}}$ , an

affine transformation and a softmax function:

$$P_{\theta}(w_{n,m+1}|w_{n,\leq m}, \mathbf{w}_{<n}, z_n) = \frac{e^{(Ow_{n,m+1})^T f_{\theta}^{\text{mlp}}(h_{n,m}^{\text{dec}})}}{\sum_{w'} e^{(Ow')^T f_{\theta}^{\text{mlp}}(h_{n,m}^{\text{dec}})}}, \quad (12)$$

where  $O \in \mathbb{R}^{|V| \times d}$  is the word embedding matrix for the output distribution with embedding dimensionality  $d \in \mathbb{N}$ .

As mentioned in the model section, the approximate posterior is conditioned on the *encoder* RNN state of the next utterance:

$$Q_{\psi}(z_n | \mathbf{w}_{\leq n}) = f_{\psi}^{\text{post}}(z_n; h_{n-1}^{\text{con}}, h_{n,M_n}^{\text{enc}}), \quad (13)$$

where  $f^{\text{post}}$  is a PDF parametrized by  $\psi$  and  $h_{n,M_n}^{\text{enc}}$  (i.e. the future state of the *encoder* RNN after processing  $\mathbf{w}_n$ ).

For the Gaussian latent variables, we use the interpolation gating mechanism described in the main text for the approximate posterior. We experimented with other mechanisms for controlling the gating variables, such as defining  $\alpha_{\mu}$  and  $\alpha_{\sigma}$  to be a linear function of the encoder. However, this did not improve performance in our preliminary experiments.

## E Appendix: Training Details

**Piecewise Constant Variable Interpolation** We conducted initial experiments with the interpolation gating mechanism for the approximate posterior of the piecewise constant latent variables. However, we found that this did not improve performance.

**Dialogue Modeling** We use the Ubuntu Dialogue Corpus v2.0 extracted January, 2016: <http://cs.mcgill.ca/~jpineau/datasets/ubuntu-corpus-1.0/>.

For the *HRED* model we found that an additional rectified linear units layer decreased performance on the validation set according to the activity F1 metric. Hence we test *HRED* without the rectified linear units layer. On the other hand, for all *VHRED* models we found that the additional rectified linear units layer improved performance on the validation set. For *P-VHRED*, we found that a final weight of one for the KL divergence terms performed best on the validation set. For *G-VHRED* and *H-VHRED*, reweighing the KL divergence terms with a final value 0.25 performed best on the validation set. We conducted preliminary experiments with  $n = 3$  and  $n = 5$  pieces,

and found that models with  $n = 3$  were easier to train. Therefore, we use  $n = 3$  pieces for both *P-VHRED* and *H-VHRED*.

For all models, we compute the log-likelihood and variational lower-bound costs starting from the second utterance in each dialogue.

## F Appendix: Additional Document Modeling Experiments

**Iterative Inference** For the document modeling experiments, our results and conclusions depend on how tight the variational lower-bound is. As such, it is in theory possible that some of our models are performing much better than reported by the variational lower-bound on the test set. Therefore, we use a non-parametric iterative inference procedure to tighten the variational lower-bound, which aims to learn a separate approximate posterior for each test example. The iterative inference procedure consists of simple stochastic gradient descent (no more than 100 steps), with a learning rate of 0.1 and the same gradient rescaling used in training. For 20 News-Groups, the iterative inference procedure is stopped on a test example if the bound does not improve over 10 iterations. For Reuters and CADE, the iterative inference procedure is stopped if the bound does not improve over 5 iterations. During iterative inference the parameters of the model, as the well as the generated prior, are all fixed. Only the gradients of the variational lower-bound with respect to generated posterior model parameters (i.e. the mean and variance of the Gaussian variables, and the piecewise components,  $a_i$ ) are used to update the posterior model for each document (using a freshly drawn sample for each inference iteration step).

Note, this form of inference is expensive and requires additional meta-parameters (e.g. a step-size and an early-stopping criterion). We remark that a simpler, and more accurate, approach to inference might perhaps be to use importance sampling.

The results based on iterative inference are reported in Table 5. As Section 6.1, we find that *H-NVDM* outperforms the *G-NVDM* model. This confirms our previous conclusions.

In our current examples, it appears that the *H-NVDM* with 5 pieces returns more general words. For example, as evidenced in Table 4, in the case of “government”, the baseline seems to value the plural form of the word (which is largely based on morphology) while the hybrid model actually

G-NVDM	H-NVDM-3	H-NVDM-5
governments	citizens	arms
citizens	rights	rights
country	governments	federal
threat	civil	country
private	freedom	policy
rights	legitimate	administration
individuals	constitution	protect
military	private	private
freedom	court	citizens
foreign	states	military

Table 4: Word query similarity test on 20 News-Groups: for the query ‘government’.

pulls out meaningful terms such as “federal”, “policy”, and “administration”.

**Approximate Posterior Analysis** We present an additional analysis of the approximate posterior on 20 News-Groups, in order to understand what the models are capturing. For a test example, we calculate the squared norm of the gradient of the KL terms w.r.t. the word embedding inputted to the approximate posterior model. The higher the squared norm of the gradients of a word is, the more influence it will have on the posterior approximation (encoder model). For every test example, we count the top 5 words with highest squared gradients separately for the multivariate Gaussian and piecewise constant latent variables.<sup>6</sup>

The results shown in Table 6, illustrate how the piecewise variables capture different aspects of the document data. The Gaussian variables were originally were sensitive to some of the words in the table. However, in the hybrid model, nearly all of the temporal words that the Gaussian variables were once more sensitive to now more strongly affect the piecewise variables, which themselves also capture all of the words that were originally missed. This shift in responsibility indicates that the piecewise constant variables are better equipped to handle certain latent factors. This effect appears to be particularly strong in the case of certain nationality-based adjectives (e.g., “american”, “israeli”, etc.). While the *G-NVDM* could model multi-modality in the data to some degree, this work would be primarily done in the model’s decoder. In the *H-NVDM*, the piecewise variables provide an explicit mechanism for capturing modes in the unknown target distribution, so it makes sense that the model would learn to use the piecewise variables instead, thus freeing up the

<sup>6</sup>Our approach is equivalent to counting the top 5 words with the highest L2 gradient norms.

Gaussian variables to capture other aspects of the data, as we found was the case with names (e.g., “jesus”, “kent”, etc.).

## G Appendix: Additional Dialogue Modeling Experiments

**Ubuntu Experiments** We present test examples — dialogue context and model responses generated using beam search — for the Ubuntu models in Table 7. The examples qualitatively illustrate the differences between models. First, we observe that *HRED* tends to generate highly generic responses compared to all the latent variable models. This supports the quantitative results reported in the main text, and suggests that modeling the latent factors through latent variables is critical for this task. Next, we observe that *H-VHRED* tends to generate relevant entities and commands — such as *mount command*, *xserver-xorg*, *static ip address* and *pulseaudio* in examples 1-4. On the other hand, *G-VHRED* tends to be better at generating appropriate verbs — such as *list*, *install*, *pastebin* and *reboot* in examples 1-3 and example 5. Qualitatively, *P-VHRED* model appears to perform somewhat worse than both *G-VHRED* and *H-VHRED*. This suggests that the Gaussian latent variables are important for the Ubuntu task, and therefore that the best performance may be obtained by combining both Gaussian and piecewise latent variables together in the *H-VHRED* model.

**Twitter Experiments** We also conducted a dialogue modeling experiment on a Twitter corpus, extracted from based on public Twitter conversations (Ritter et al., 2011). The dataset is split into training, validation, and test sets, containing respectively 749,060, 93,633 and 9,399 dialogues each. On average, each dialogue contains about 6 utterances (dialogue turns) and about 94 words. We pre-processed the tweets using byte-pair encoding (Sennrich et al., 2016) with a vocabulary consisting of 5000 sub-words.

We trained our models with a learning rate of 0.0002 and mini-batches of size 40 or 80.<sup>7</sup> As for the Ubuntu experiments, we used a variant of truncated back-propagation and apply gradient clipping. We experiment with *G-VHRED* and *H-VHRED*. Similar to (Serban et al., 2017b), we use a bidirectional GRU RNN *encoder*, where the forward and backward RNNs each have 1000 hid-

<sup>7</sup>We had to vary the mini-batch size to make the training fit on GPU architectures with low memory.

20-NG	Sampled	SGD-Inf	RCV1	Sampled	SGD-Inf
<i>LDA</i>	1058	---	<i>G-NVDM</i>	905	837
<i>RSM</i>	953	---	<i>H-NVDM-3</i>	865	807
<i>docNADE</i>	896	---	<i>H-NVDM-5</i>	<b>833</b>	<b>781</b>
<i>SBN</i>	909	---			
<i>fDARN</i>	917	---	<b>CADE</b>	<b>Sampled</b>	<b>SGD-Inf</b>
<i>NVDM</i>	836	---	<i>G-NVDM</i>	339	230
<i>G-NVDM</i>	651	588	<i>H-NVDM-3</i>	<b>258</b>	<b>193</b>
<i>H-NVDM-3</i>	607	546	<i>H-NVDM-5</i>	294	209
<i>H-NVDM-5</i>	<b>566</b>	<b>496</b>			

Table 5: Comparative test perplexities on various document datasets (50 latent variables). Note that document probabilities were calculated using 10 samples to estimate the variational lower-bound.

den units. We experiment with *context* RNN encoders with 500 and 1000 hidden units, and find that that 1000 hidden units reach better performance w.r.t. the variational lower-bound on the validation set. The *encoder* and *context* RNNs use layer normalization (Ba et al., 2016). We experiment with *decoder* RNNs with 1000, 2000 and 4000 hidden units (LSTM cells), and find that 2000 hidden units reach better performance. For the *G-VHRED* model, we experiment with latent multivariate Gaussian variables with 100 and 300 dimensions, and find that 100 dimensions reach better performance. For the *H-VHRED* model, we experiment with latent multivariate Gaussian and piecewise constant variables each with 100 and 300 dimensions, and find that 100 dimensions reach better performance. We drop words in the decoder with a fixed drop rate of 25% and multiply the KL terms in the variational lower-bound by a scalar, which starts at zero and linearly increases to 1 over the first 60,000 training batches. Note, unlike the Ubuntu experiments, the final weight of the KL divergence is exactly one (hence the bound is tight).

Our hypothesis is that the piecewise constant latent variables are able to capture multi-modal aspects of the dialogue. Therefore, we evaluate the models by analyzing what information they have learned to represent in the latent variables. For each test dialogue with  $n$  utterances, we condition each model on the first  $n - 1$  utterances and compute the latent posterior distributions using all  $n$  utterances. We then compute the gradients of the KL terms of the multivariate Gaussian and piecewise constant latent variables w.r.t. each word in the dialogue. Since the words vectors are discrete, we compute the sum of the squared gradi-

ents w.r.t. each word embedding. The higher the sum of the squared gradients of a word is, the more influence it will have on the posterior approximation (encoder model). For every test dialogue, we count the top 5 words with highest squared gradients separately for the multivariate Gaussian and piecewise constant latent variables.<sup>8</sup>

The results are shown in Table 8. The piecewise constant latent variables clearly capture different aspects of the dialogue compared to the Gaussian latent variables. The piecewise constant variable approximate posterior encodes words related to time (e.g. weekdays and times of day) and events (e.g. parties, concerts, Easter). On the other hand, the Gaussian variable approximate posterior encodes words related to sentiment (e.g. laughter and appreciation) and acronyms, punctuation marks and emoticons (i.e. smilies). We also conduct a similar analysis on the document models evaluated in Sub-section 6.1, the results of which may be found in the Appendix.

<sup>8</sup>Our approach is equivalent to counting the top 5 words with the highest L2 gradient norms. We also did some experiments using L1 gradient norms, which showed similar patterns.

Word	G-NVDM	H-NVDM-5		Word	G-NVDM	H-NVDM-5	
<b>Time-related</b>	<b>G-KL</b>	<b>G-KL</b>	<b>P-KL</b>	<b>Names</b>	<b>G-KL</b>	<b>G-KL</b>	<b>P-KL</b>
months	23	33	<b>40</b>	henry	33	<b>47</b>	39
day	28	32	<b>35</b>	tim	<b>32</b>	27	11
time	<b>55</b>	22	<b>40</b>	mary	26	<b>51</b>	30
century	<b>28</b>	13	<b>19</b>	james	40	<b>72</b>	30
past	<b>30</b>	18	<b>28</b>	jesus	28	<b>87</b>	39
days	<b>37</b>	14	<b>19</b>	george	26	<b>56</b>	29
ahead	<b>33</b>	20	<b>33</b>	keith	65	<b>94</b>	61
years	<b>44</b>	16	<b>38</b>	kent	51	<b>56</b>	15
today	46	27	<b>71</b>	chris	38	<b>55</b>	28
back	31	30	<b>47</b>	thomas	19	<b>35</b>	19
future	<b>20</b>	15	<b>20</b>	hitler	10	<b>14</b>	9
order	<b>42</b>	14	<b>26</b>	paul	25	<b>52</b>	18
minute	15	34	<b>40</b>	mike	38	<b>76</b>	40
began	<b>16</b>	5	<b>13</b>	bush	<b>21</b>	20	14
night	<b>49</b>	12	<b>18</b>				
hour	<b>18</b>	<b>17</b>	16	<b>Adjectives</b>	<b>G-KL</b>	<b>G-KL</b>	<b>P-KL</b>
early	42	42	<b>69</b>	american	<b>50</b>	12	<b>40</b>
yesterday	25	26	<b>36</b>	german	<b>25</b>	21	<b>22</b>
year	<b>60</b>	17	<b>21</b>	european	20	17	<b>27</b>
week	28	54	<b>58</b>	muslim	19	7	<b>23</b>
hours	20	26	<b>31</b>	french	11	<b>17</b>	<b>17</b>
minutes	<b>40</b>	34	<b>38</b>	canadian	<b>18</b>	10	<b>16</b>
months	23	33	<b>40</b>	japanese	16	9	<b>24</b>
history	<b>32</b>	18	<b>28</b>	jewish	<b>56</b>	37	<b>54</b>
late	41	<b>45</b>	31	english	19	16	<b>26</b>
moment	<b>23</b>	<b>17</b>	16	islamic	14	18	<b>28</b>
season	<b>45</b>	29	<b>37</b>	israeli	<b>24</b>	14	<b>18</b>
summer	29	28	<b>31</b>	british	<b>35</b>	15	<b>17</b>
start	30	14	<b>38</b>	russian	14	19	<b>20</b>
continue	21	32	<b>34</b>				
happened	22	27	<b>35</b>				

Table 6: Approximate posterior word encodings (20-NG). For P-KL, we bold every case where piecewise variables showed greater word sensitivity than Gaussian variables w/in the same hybrid model.

Dialogue Context (History)	Response
<p>Hi . I am installing ubuntu now in my new laptop . In " something else " partitioning , what mount point should I set for a drive which is not root or not home ... → It 's up to you , just choose a directory that will remind you of the contents of that partition . E.G. : if it 's the Windows partition , use /windows . → it 's a new harddrive with full free space . I bought it without windows preinstalled . I want to create drives in which I will only store files .. I mean , not root or not home . What mount point do I set for it ? " /mount " is not shown in drop down menu sorry . I mean /mount I mean , in my desktop , extra drives are mounted in /media do you understand my problem ? Sorry , English is not my native language . → I do :) Just use a directory with a meaningful name and preferably in the root , for example " /files " . Choose something like " /files " . → / for root , /home for home , but what mount point for a file storage partition . wait , let me see if there is a " /files " , thank you .</p>	<p><b>HRED:</b> No problem .  <b>G-VHRED:</b> It should be in the list of the drive . If you want to be able to mount it .  <b>P-VHRED:</b> If you want to mount it , you can use the mount command .  <b>H-VHRED:</b> As far as I know , you can use the " mount " command to mount the partitions .</p>
<p>I have a problem installing Ubuntu 12.04 desktop . Tried desktop with USB -&gt; hangs . alternate with usb -&gt; hangs . alternate with cd -&gt; hangs ( even faster ) . I tried using different usb drives . Could someone please help me ? → how are you making the usb ? and whats the pc specs . → I used Universal usb installer → what exactly does it do when it ' hangs ' → Specs : CPU core i7 @ 3.4 Ghz , Sapphire Radeon HD 7800 , 8GB RAM → ATI video ? tried the NOMODESET option yet ? → nothing at all , I also can't witch to console using Ctrl + alt + F not yet Will try that , thanks</p>	<p><b>HRED:</b> you could try the alternate install cd .  <b>G-VHRED:</b> it might be a problem with your graphics card . You could also try the alternate install cd and see if it works .  <b>P-VHRED:</b> try ' nomodeset ' or ' nomodeset '  <b>H-VHRED:</b> only thing you can do is run the command ' sudo dpkg-reconfigure xserver-xorg ' and see if you get any errors .</p>
<p>anyone know how to add a route in ubuntu ? getting an error on adding a route , and I'm not sure why → pastebin your command and problem → &lt;url &gt;</p>	<p><b>HRED:</b> can you pastebin the output of " route -n " please  <b>G-VHRED:</b> if you want to use pastebin to pastebin the output of " sudo route -n " please pastebin the output of " route "  <b>P-VHRED:</b> what is the output of : route -n  <b>H-VHRED:</b> only if you have a static ip address , you need to specify the ip address of the gateway</p>
<p>Any idea why empathy 's not playing notification sounds ? Even though I have ' em ticked in preferences → restarted it yet ? → yar → check pulse to see if the application is muted for some reason ? well Sound settings . → Had sound effects turned off in sound settings , didn't realize that controlled other applications → Ah yea , ive done it a few time it 's annoying My favorite though is recently pulse has been freezing on my desktop and audio will just not be adjustable for like ... 30 seconds or so</p>	<p><b>HRED:</b> I don't think it 's a problem with pulseaudio , but I'm not sure how to fix it  <b>G-VHRED:</b> Thank you for your help . I will try that . Thanks for the suggestion .  <b>P-VHRED:</b> I think it 's a bug with pulseaudio , but I don't know how to fix it .  <b>H-VHRED:</b> as far as I can tell , I have a problem with pulseaudio , but I don't know how to fix it .</p>
<p>Question about setting up drives : Running 12.04 LTS . I have two drives . An SSD that holds the OS for booting , and a second drive " data " for all the data , home directories , etc . The data drive is encrypted with LUKS . With my current setup , whenever I reboot , I get a password entry screen to decrypt the data drive , where my home directory is . This works fine . The problem is that if I reboot remotely , I'm stuck as I can't enter password . If I keep my home directory on the SSD , there won't be enough room . If I use a symbolic link , and the encrypted drive isn't mounted , that creates errors and issues . Suggestions ?? → yes , you can redirect the → Message was cut off ...</p>	<p><b>HRED:</b> if you want to see the data on the drive , you can mount it in /etc/fstab  <b>G-VHRED:</b> it should be a good place to do that , but if you want to do a clean install , you will need to reboot  <b>P-VHRED:</b> it 's not a problem , it 's just a bug  <b>H-VHRED:</b> if you want to see the data on the drive , you can mount it in /etc/fstab</p>
<p>hey guys , does anyone know why i cannot connect to ssh via my external IP ? (connection refused ) but ssh -v localhost works fine ? → Have you forwarded the port on your router to the machine you're trying to access ? → It is a work machine , so am not sure of the rules</p>	<p><b>HRED:</b> You need to forward port 22 to your router to forward port 22 to the server .  <b>G-VHRED:</b> That 's odd . What are you trying to do ? Can you pastebin the output of " sudo netstat " to " pastebin " please ?  <b>P-VHRED:</b> Can you pastebin the output of " sudo apt-get install openssh-server "?  <b>H-VHRED:</b> Even if it 's not working , then you need to set the port forward to your router .</p>

Table 7: Ubuntu model examples. The → token indicates a change of turn.

Word Time-related	G-VHRED			Word Event-related	H-VHRED		
	G-KL	G-KL	P-KL		G-KL	G-KL	P-KL
monday	3	5	<b>10</b>	school	9	16	<b>50</b>
tuesday	2	3	<b>7</b>	class	11	16	<b>27</b>
wednesday	4	11	<b>13</b>	game	20	26	<b>41</b>
thursday	2	3	<b>9</b>	movie	12	20	<b>41</b>
friday	9	18	<b>26</b>	club	13	22	<b>28</b>
saturday	6	6	<b>13</b>	party	8	10	<b>32</b>
sunday	2	2	<b>9</b>	wedding	7	13	<b>23</b>
weekend	8	16	<b>32</b>	birthday	12	20	<b>23</b>
today	18	28	<b>56</b>	easter	15	15	<b>23</b>
night	16	31	<b>68</b>	concert	7	16	<b>20</b>
tonight	32	36	<b>47</b>	dance	11	12	<b>21</b>

  

Word Sentiment -related	G-VHRED			Word Acronyms, Punctuation Marks & Emoticons	H-VHRED		
	G-KL	G-KL	P-KL		G-KL	G-KL	P-KL
good	<b>72</b>	<b>73</b>	44	lol	<b>394</b>	<b>358</b>	312
love	<b>102</b>	<b>101</b>	38	omg	<b>52</b>	<b>45</b>	19
awesome	<b>26</b>	<b>44</b>	39	.	386	558	<b>1009</b>
cool	14	28	<b>29</b>	!	<b>648</b>	<b>951</b>	525
haha	<b>132</b>	<b>101</b>	75	?	<b>507</b>	<b>851</b>	221
hahaha	<b>60</b>	<b>48</b>	24	*	<b>108</b>	<b>54</b>	19
amazing	<b>14</b>	<b>38</b>	33	xd	<b>28</b>	<b>42</b>	26
thank	<b>137</b>	<b>153</b>	29	♡	<b>56</b>	<b>42</b>	24

Table 8: Approximate posterior word encoding on Twitter. The numbers are computed by counting the number of times each word is among the 5 words with the largest sum of squared gradients of the Gaussian KL divergence (G-KL) and piecewise constant KL divergence (P-KL)