
Supporting Content Marketing with Natural Language Generation

WEB APPENDIX

6/28/2021 12:06:30 AM

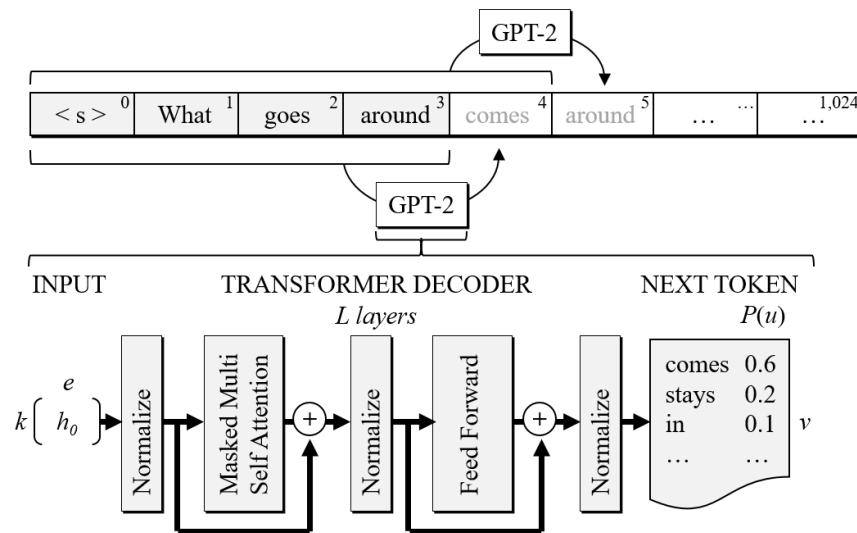
Index

1	Technical Modeling and Validation Notes	3
1.1	GPT-2 Model Description.....	3
1.2	Description of Essential Software Features	5
1.3	Applied Uniqueness, Naturality & Readability Measures	7
1.4	External Validation of Method Assumptions and Quality Score.....	8
1.5	Validation of Method Fine-Tuning Process	13
1.6	External Validation of Method Performance	16
1.7	Alternative Quality Score Weighting Performance Assessment.....	17
1.8	Machine Generated Sample Content.....	19
2	IT Service Industry Field Study	23
2.1	Keywords Used & Keyword Statistics.....	23
2.2	Experimental Setup	24
2.3	Performance of the Experimental Groups.....	28
2.3.1	Quality Scores of Experimental Groups' Content	28
2.3.2	Post Hoc Test for Achieved Search Engine Rankings.....	31
2.3.3	Content Ranking for Sub-Keywords Assessment.....	32
2.3.4	Additional Keywords Performance.....	33
2.4	Providing Quality Score Feedback to Revise Content	36
2.5	Consumers' Content Perceptions	39
2.6	Website Engagement.....	44
2.7	Content Production Cost Calculation Details.....	46
3	Education Sector Field Study	48
3.1	Keywords to Optimize for & Keyword Statistics	48
3.2	Supplemental Content Performance Tests	51
4	Keyword Boundary Conditions.....	54
5	Appendix References	56

1 Technical Modeling and Validation Notes

1.1 GPT-2 Model Description

To provide the intuition behind transformer-based NLG models, we briefly illustrate the mechanics of the popular GPT-2 model. Given a sequence of tokens with context window size k , $U=(u_{-k}, \dots, u_{-1})$, the objective of the autoregressive model GPT-2 is to accurately “predict” the next likely word¹ (Figure W1.1) by sampling from a probability distribution over its entire learned vocabulary (consisting of 50,257 tokens) conditional on the given word sequence and on a pre-trained neural network with parameters Θ . Model pre-training tries to maximize the likelihood in equation (W1) for an unsupervised corpus of words (\mathcal{U}) (Radford et al. 2018).

Figure W1.1: The GPT-2 Model²

¹ For ease of discussion, we describe the model in terms of “words.” GPT-2 is derived using BPE (Byte Pair Encoding) and tokens (i.e., learned and encoded pieces of words).

² Visualization derived from Radford et al. (2018), and adapted to depict the updated GPT-2 architecture.

$$L_1(\mathcal{U}) = \sum_i \log P(u_i | u_{i-k}, \dots, u_{i-1}; \Theta) \quad (\text{W1})$$

$$h_0 = UW_e + W_p \quad (\text{W2})$$

$$h_i = \text{transformer_block}_{(h_{i-1})} \forall i \in [1, L] \quad (\text{W3})$$

$$P(u) = \text{softmax}(h_L W_e^T) \quad (\text{W4})$$

In essence, GPT-2 relies on word and given context meaning information to generate its output distribution over its vocabulary. More specifically, the data input consists of a matrix h_0 (W2), where the given word sequence U , word meaning information in terms of word embeddings W_e , and sequential word position information in terms of position embeddings W_p are combined. As illustrated in Figure W1.1, information from h_0 is extracted, transformed, added and normalized multiple times (to ease processing), and projected into the embedding space e by L layers of decoder transformer blocks (W3). This information includes the extent of putting attention on a given word sequence using multi-headed self attention (Vaswani et al. 2017), and high dimensional hidden language states to shift the focus in the embedding space e to recreate natural word sequences from position-wise feed forward neural networks. The output of the final block h_L projects all this information into the embedding space and is multiplied with GPT-2's original (unconditional) transposed word embeddings matrix W_e^T to assess which word from the GPT-2 vocabulary best matches the information contained in h_L (W4). The multiplication of h_L and W_e^T can be thought of as a similarity or matching between the embedding space distribution of the output of h_L (containing meaning, position, attention, and hidden language states information) and the unconditional embedding space distribution of each respective vocabulary word. More similarity of a vocabulary word in terms of its embedding to h_L will result in a higher probability in GPT-2's output distribution. GPT-2 then obtains a

probability distribution over its vocabulary $P(u)$ (W4) and can sample the upcoming word in the sequence from the most likely words in $P(u)$.

Using the above procedure, GPT-2 learned and stored word probabilities for given word sequences represented in its 345 million parameters (including embeddings, attention weight matrices and Θ) using 8 million English text documents with a broad topical variety. Neural network parameters Θ were first initialized and then trained on batches of 512 sequences. The loss function refers to the language modeling cross entropy loss, where 1 is assigned to the word that appears next (u_i) in the training sequence (e.g., “comes” in Figure W1.1), and 0 to all other words in GPT-2’s vocabulary, and compare the log transformed GPT-2 softmaxed output probability value P_u for that respective word to appear next. A loss of 0 means the GPT-2 prediction was in perfect accordance with the actual next word (i.e., 1), the higher the deviation of the GPT-2 prediction (e.g., $1-0.6 = 0.4$ for “comes” in Figure W1.1) to the actual word, the more the loss value increases. During training, GPT-2 performs this process on batches and mini-batches of several sequences before updating Θ .

1.2 Description of Essential Software Features

For our content generation method to work “on the fly,” a number of features have been incorporated at each stage of our algorithm for it to work automatically and reliably for any specified keyword. We summarize the most essential high-level features in Table W1.2.

Table W1.2: Developed Content Writing Machine Software Features

Step	Feature	Feature description
Ranking & links crawling	Crawler status updates	The crawler provides live-status updates
	Human behavior simulation	The crawler simulates human behavior to not get blocked by the search engine
	Organic links detection	Organic links are automatically detected and other links (e.g., paid ads) are discarded
	Duplicate entries detection	Duplicate ranking entries (e.g., featured content snippet) are detected and discarded
	Link correction (abbreviations, prefixes, etc.)	Abbreviated organic link prefixes are detected and corrected
	Data handling & saving	The crawler handles and saves the data in a structured way
Content scraping	Scraper status updates	The scraper provides live-status updates
	Local user client simulation	The scraper simulates an actual user client including a local OS, a local browser and a record of cookies
	Enhanced SSL protocols	The scraper uses enhanced SSL-protocols
	Website main content recognition	The scraper analyzes the HTML-files as well as common text patterns to find the main content on the webpage (i.e., discarding content of footers, main menus, etc.)
	Code cleanup	The scraper detects HTML, CSS & Java script code and cleans the main content from it
	Text cleanup	The scraper detects unwanted text snippets and patterns (i.e., big empty spaces, unintended line breaks, citing, remarks, etc.) and cleans the main content from it
	Server failure messages cleanup	The scraper universally detects server messages (e.g., “you are not allowed to crawl this website”, “502 server error”, etc.) and discards these
	Multi-redundancy	The scraper consists of multiple safety lines for full automation including error detection and handling, timeouts on jobs, etc.
	Auto-output	The scraper automatically outputs a txt-file that includes all information for retraining, including the main content, the targeted main keyword, and automatically generated special tags
Fine-tuning & content generation	Dynamic retraining & content generation	Our method performs a dynamic fine-tuning (i.e., several retraining, model checkpoint-saving and text generation steps where in each step, the retraining is continued to cover the full spectrum of model fine-tuning and fitting on the data)
	Fallback model	Our method uses the base model as fallback (i.e., if the retraining material is corrupt, it still generates and provides well written texts from an early retraining phase)
	H1, auto seed-word, and auto tagging	Our method takes the specified main keyword automatically as text headline and as seed-word for text generation for increased text consistency and topic focus
	Model checkpoints & file saving	Our method automatically performs several model-checkpoint savings and outputs the generated content in structured txt-files
Content selection & output	Auto data handling	The generated content is automatically cleaned and handled
	Auto quality score calculation	The quality score for text selection is automatically calculated
	Redundancy (error handling)	Errors are detected, outputted and appropriately handled
	Intuitive ordered list output for humans	An intuitive annotated output in the form of an ordered list of suggested generated and selected content is provided to a human reviewer

1.3 Applied Uniqueness, Naturality & Readability Measures

Without loss of generality, the quality score we present in the article could be adapted to incorporate other linguistic components. The software tool employed in our empirical application studies implements the components content uniqueness (s_d), naturality similarity (s_n) and readability similarity (s_r) as follows:

Uniqueness measurement (s_d). For our quality score (qs_g), we derive a uniqueness measure (s_d) to assess if the content is sufficiently unique for the search engine. In addition to the definitions around formula (3) in the main manuscript, we apply a critical value (s_{cv}) to ensure that the generated content is sufficiently unique based on the length of the keyword (kw) and parameter b .

$$s_{cv} = (100 - (100/(kw + 1)^b))/100 \quad (W5)$$

By implementing this non-compensatory filtering rule we ascertain that content that fails to achieve this minimum level of uniqueness is discarded from further content selection. The value b determines the factor of increasing conservativeness the larger the n -gram size ($kw+1$), as repeating small sized n -grams is less of a concern than repeating large sized n -grams (W5). In our setup, we set b to 1.1 after an evaluation phase in which we look at a) the machine output, b) acceptable duplicate rates in human content impressions, and c) content retaining rates for the whole range of common n -gram sizes. For example, that means that with an n -gram size of 3, $s_{cv} \sim .70$ (i.e., 70% unique), an n -gram size of 5, $s_{cv} \sim .82$ (i.e., 82% unique), and an n -gram-size of 7, $s_{cv} \sim .88$ (i.e., 88% unique).

Naturality similarity measures (s_n). To quantify the naturality similarity between the generated content and the top ranked search results, we applied 12 linguistic measures which assess the lexical richness and composition of a text using the R package [languageR](#).

Specifically, we use the following measures: tokens, types, hapax legomena, dis legomena, tris legomena, Yule's K, Zipf's R, Type-Token-Ratio, Herdan's C, Guiraud's R, Sichel's S, Lognormal. More information on the precise meaning, practical examples and literature sources can be found in Baayen and Shafaei-Bajestan (2019).

Readability similarity measures (s_r). For the readability similarity measure, we applied 46 pre-existing measures of readability contained in the R package [quanteda](#) (see Benoit et al. 2018). We make use of the following measures: ARI, Bormuth.MC, Bormuth.GP, Coleman, Coleman.C2, Coleman.Liau.ECP, Dale.Chall, Dale.Chall.PSK, Danielson.Bryan, Dickes.Steiwer, DRP, ELF, Farr.Jenkins.Paterson, Flesch, Flesch.PSK, Flesch.Kincaid, FOG, FOG.PSK, FOG.NRI, FORCAST, FORCAST.RGL, Fucks, Linsear.Write, nWS, nWS.2, nWS.3, nWS.4, RIX, Scrabble, SMOG, SMOG.C, Spache, Spache.old, Strain, Traenkle.Bailer, W, St, C, Sy, W3Sy, W2Sy, W_1Sy, W6C, W7C, Wlt3Sy, W_wl.Dale.Chall. More information on the precise meaning, calculation and literature sources can be found in Benoit et al. (2018).

1.4 External Validation of Method Assumptions and Quality Score

Before using our method in a field application, we empirically test and confirm that the highest-ranking websites in the search engine indeed score highest in terms of our developed quality score components. For this task, we used around 8,500 relevant keywords and about 1.42 million ranked websites from 4 main industry sectors and 36 specific industries: details on the distribution of these keywords across industries and the number of scraped website content are reported in Table W1.4.1, columns (1) and (2). Using Wilcoxon rank sum group comparison tests, Table W1.4.2 illustrates that the poorer the search engine ranking, the lower the quality

scores compared to the top 10 ranked content tends to be for all quality score components with the exception for content uniqueness (s_d).

Recall that the content uniqueness score s_d compares a given piece of content to the top 10 search results. The observed pattern in Table W1.4.2 suggests that a given result from outside of the top 10 set is more unique compared to the top 10 set than a result from the top 10 set is unique compared to other top 10 results. As noted in the manuscript, this may arise because the top 10 ranked websites consistently reflect similar topics. A contributor to the similarity among the top 10 ranked websites are the topics and keywords discussed in the content, as reflected by the higher scores for (s_a and s_k , respectively) compared to lower ranked websites. In assessing uniqueness relative to the content of the top 10 ranked websites, as these sites tend to cover similar topics and keywords compared to lower ranked sites, it is not surprising to find that the uniqueness score among the top 10 ranked websites is lower than the uniqueness score associated with lower ranked websites. Thus, we can ascertain that fine-tuning on the top 10 ranked websites' content will produce the most optimal content, and supporting the use of our quality score as a measure of content optimality. We aggregate the search engine results into groups (i.e., top 10, search engine ranks 11-20, search engine ranks 21-99, search engine ranks 100-200) to summarize the results.

Table W1.4.1: Empirical Setup for Validating Method and Quality Score Assumptions

Industry Sector	Industry	(1) Number of Keywords	(2) Number of Scraped Rankings & Websites	(3) Number of Selected Keywords	(4) Number of Generated Texts
I.	Coal Mining	100	14,678	5	5,000
	Forestry	501	87,537	9	9,000
	Grazing	100	18,021	10	10,000
	Hunting	100	17,621	7	7,000
	Fishing	500	77,210	10	10,000
	Quarrying	176	18,448	8	8,000
II.	Automobile production	270	42,303	10	10,000
	Textile production	150	26,960	9	9,000
	Chemical engineering	230	43,288	8	8,000
	Aerospace production	250	57,149	10	10,000
	Energy utilities	150	29,767	10	10,000
	Breweries & bottlers	150	30,691	9	9,000
	Construction	150	21,757	7	7,000
	Ship building	70	14,058	9	9,000
	Jewelries	245	45,097	9	9,000
III.	Retailing	150	27,717	9	9,000
	Transportation	450	60,222	9	9,000
	Restaurants	230	32,539	9	9,000
	Clerical service	300	49,188	9	9,000
	Mass media	300	39,784	9	9,000
	Tourism	300	41,174	10	10,000
	Insurance	150	27,581	10	10,000
	Banking	270	44,007	9	9,000
	Healthcare	150	30,478	10	10,000
	Law	230	43,717	9	9,000
	IT service	324	50,670	19	19,000
	Art & galleries	150	27,167	9	9,000
	Cafes	230	35,382	9	9,000
	Grocery stores	500	80,814	10	10,000
	Media agencies	150	29,180	10	10,000
IV.	Government	300	50,074	9	9,000
	University	349	54,775	11	11,000
	Culture	300	57,704	9	9,000
	Libraries	100	15,715	9	9,000
	Research	100	9,938	10	10,000
	Education	278	62,518	10	10,000

Table W1.4.2: External Validation of Method Assumptions Statistics

Industry Sector	Ranks of Content Compared to Top 10	Topic (s_a) ¹	Keywords (s_k) ¹	Uniqueness (s_d) ¹	Readability similarity (s_r) ¹	Naturality similarity (s_n) ¹
I.	Top 10	.27 (.16)	.23 (.23)	.93 (.22)	.74 (.57)	.75 (.50)
	11 - 20	.23 (.17)**	.18 (.23)**	.96 (.11)**	.70 (.62)**	.58 (.58)**
	21 - 99	.18 (.15)**	.13 (.20)**	.96 (.09)**	.65 (.55)**	.58 (.58)**
	100 - 200	.15 (.15)**	.09 (.20)**	.97 (.09)**	.70 (.62)**	.67 (.58)**
II.	Top 10	.31 (.17)	.26 (.22)	.95 (.15)	.70 (.57)	.67 (.50)
	11 - 20	.25 (.16)**	.20 (.22)**	.97 (.09)**	.62 (.62)**	.58 (.50)**
	21 - 99	.22 (.17)**	.16 (.21)**	.97 (.08)**	.59 (.59)**	.58 (.50)**
	100 - 200	.17 (.15)**	.11 (.21)**	.97 (.07)**	.57 (.57)**	.50 (.50)**
III.	Top 10	.35 (.22)	.31 (.30)	.94 (.17)	.72 (.60)	.75 (.50)
	11 - 20	.29 (.21)**	.25 (.29)**	.96 (.10)**	.70 (.60)**	.67 (.58)**
	21 - 99	.23 (.20)**	.17 (.26)**	.97 (.08)**	.64 (.60)**	.58 (.58)**
	100 - 200	.18 (.17)**	.10 (.22)**	.98 (.06)**	.57 (.62)**	.50 (.58)**
IV.	Top 10	.31 (.20)	.26 (.27)	.95 (.10)	.72 (.60)	.62 (.58)
	11 - 20	.27 (.20)**	.21 (.25)**	.97 (.08)**	.68 (.57)**	.58 (.58)**
	21 - 99	.22 (.19)**	.14 (.21)**	.97 (.07)**	.62 (.60)**	.57 (.58)**
	100 - 200	.16 (.16)**	.07 (.18)**	.97 (.06)**	.62 (.59)**	.42 (.67)**

¹Reported numbers are group medians and IQRs in parentheses. Statistical significance codes come from Wilcoxon rank sum 2-group comparison tests between top 10 ranked websites and the content with specific rankings as stated in column 2; statistical significance codes (one-tailed): *0.05 level, **0.01 level; assumptions (e.g., non-normality of data) for all Wilcoxon rank-sum 2-group comparison tests are confirmed.

The results of Table W1.4.2 are consistent for smaller sets of search engine rankings (which correspond to a given page of search engine results) for the single exemplary industry sector III (Table W1.4.3) and for specific industries (e.g., tourism) (Table W1.4.4). The observed patterns are consistent across industries and the tests for significance are insensitive to industry sector aggregations and/or data-groupings; more detailed industry-specific results are available from the authors upon request.

Table W1.4.3: External Validation of Method Assumptions Statistics for Industry Sector III

Industry Sector	Ranks of Content Compared to Top 10	Topic (s_a) ¹	Keywords (s_k) ¹	Uniqueness (s_d) ¹	Readability similarity (s_r) ¹	Naturality similarity (s_n) ¹
III.	Top 10	.35 (.22)	.31 (.30)	.94 (.17)	.72 (.60)	.75 (.50)
	11 - 20	.29 (.21)**	.25 (.29)**	.96 (.10)**	.70 (.60)**	.67 (.58)**
	21 - 30	.27 (.20)**	.22 (.27)**	.96 (.09)**	.66 (.60)**	.58 (.50)**
	31 - 40	.25 (.20)**	.20 (.27)**	.97 (.09)**	.66 (.62)**	.58 (.58)**
	41 - 50	.24 (.20)**	.19 (.26)**	.97 (.08)**	.63 (.59)**	.58 (.58)**
	51 - 60	.23 (.19)**	.17 (.25)**	.97 (.08)**	.64 (.57)**	.58 (.58)**
	61 - 70	.23 (.19)**	.16 (.25)**	.97 (.08)**	.63 (.62)**	.58 (.67)**
	71 - 80	.22 (.19)**	.15 (.25)**	.97 (.08)**	.66 (.62)**	.58 (.67)**
	81 - 90	.21 (.18)**	.15 (.24)**	.97 (.08)**	.64 (.64)**	.50 (.58)**
	91 - 100	.19 (.18)**	.12 (.23)**	.97 (.08)**	.62 (.62)**	.50 (.58)**
	101 - 110	.19 (.19)**	.12 (.24)**	.97 (.07)**	.61 (.62)**	.50 (.58)**
	111 - 120	.19 (.18)**	.11 (.23)**	.98 (.07)**	.62 (.59)**	.50 (.58)**
	121 - 130	.20 (.18)**	.12 (.24)**	.98 (.07)**	.57 (.55)**	.50 (.58)**
	131 - 140	.18 (.19)**	.11 (.25)**	.97 (.08)**	.60 (.55)**	.50 (.58)**
	141 - 150	.17 (.16)**	.10 (.20)**	.97 (.07)**	.62 (.59)**	.50 (.58)**
	151 - 160	.16 (.16)**	.09 (.20)**	.97 (.07)**	.53 (.57)**	.42 (.67)**
	161 - 170	.16 (.16)**	.09 (.21)**	.97 (.08)**	.57 (.62)**	.50 (.58)**
	171 - 180	.16 (.15)**	.08 (.19)**	.97 (.07)**	.55 (.56)**	.50 (.58)**
	181 - 190	.16 (.15)**	.08 (.19)**	.97 (.07)**	.61 (.62)**	.58 (.58)**
	191 - 200	.15 (.14)**	.07 (.17)**	.97 (.07)**	.49 (.55)**	.50 (.58)**

¹Reported numbers are group medians and IQRs in parentheses. Statistical significance codes come from Wilcoxon rank sum 2-group comparison tests between top 10 ranked websites and the content with specific rankings as stated in column 2; statistical significance codes (one-tailed): *0.05 level, **0.01 level; assumptions (e.g., non-normality of data) for all Wilcoxon rank-sum 2-group comparison tests are confirmed.

Table W1.4.4: External Validation of Method Assumptions for the Tourism Sector

Industry	Ranks of Content Compared to Top 10	Topic (s_a) ¹	Keywords (s_k) ¹	Uniqueness (s_d) ¹	Readability similarity (s_r) ¹	Naturality similarity (s_n) ¹
Tourism	Top 10	.37 (.32)	.36 (.40)	.88 (.27)	.63 (.60)	.58 (.58)
	11 - 20	.32 (.21)**	.30 (.29)*	.95 (.18)**	.60 (.53)	.50 (.63)*
	21 - 30	.30 (.22)**	.26 (.28)**	.95 (.09)**	.48 (.54)**	.42 (.58)**
	31 - 40	.28 (.32)**	.26 (.26)**	.95 (.11)**	.42 (.48)**	.42 (.50)**
	41 - 50	.28 (.20)**	.26 (.27)**	.96 (.08)**	.49 (.53)**	.41 (.50)**
	51 - 60	.25 (.16)**	.21 (.24)**	.97 (.08)**	.53 (.51)**	.42 (.58)**
	61 - 70	.25 (.19)**	.23 (.27)**	.95 (.11)**	.43 (.52)**	.42 (.50)**
	71 - 80	.23 (.18)**	.18 (.27)**	.97 (.08)**	.43 (.53)**	.33 (.50)**
	81 - 90	.20 (.21)**	.14 (.25)**	.97 (.07)**	.51 (.55)**	.33 (.50)**
	91 - 100	.23 (.21)**	.17 (.28)**	.96 (.09)**	.43 (.59)**	.33 (.48)**
	101 - 110	.19 (.19)**	.15 (.26)**	.97 (.10)**	.43 (.57)**	.33 (.58)**
	111 - 120	.16 (.17)**	.10 (.22)**	.97 (.08)**	.34 (.52)**	.25 (.35)**
	121 - 130	.16 (.18)**	.07 (.22)**	.96 (.11)**	.32 (.51)**	.33 (.46)**
	131 - 140	.14 (.15)**	.08 (.20)**	.97 (.07)**	.36 (.49)**	.42 (.46)**
	141 - 150	.17 (.19)**	.13 (.23)**	.97 (.10)**	.33 (.55)**	.29 (.56)**
	151 - 160	.12 (.10)**	.06 (.14)**	.97 (.07)**	.46 (.40)**	.25 (.25)**
	161 - 170	.10 (.11)**	.04 (.10)**	.96 (.09)**	.33 (.45)**	.38 (.50)**
	171 - 180	.16 (.12)**	.10 (.16)**	.97 (.11)**	.51 (.49)**	.33 (.48)**
	181 - 190	.12 (.19)**	.04 (.18)**	.97 (.03)**	.46 (.62)**	.25 (.50)**
	191 - 200	.13 (.11)**	.10 (.14)**	.98 (.07)**	.59 (.48)	.50 (.58)

¹Reported numbers are group medians and IQRs in parentheses. Statistical significance codes come from Wilcoxon rank sum 2-group comparison tests between top 10 ranked websites and the content with specific rankings as stated in column 2; statistical significance codes (one-tailed): *0.05 level, **0.01 level; assumptions (e.g., non-normality of data) for all Wilcoxon rank-sum 2-group comparison tests are confirmed.

1.5 Validation of Method Fine-Tuning Process

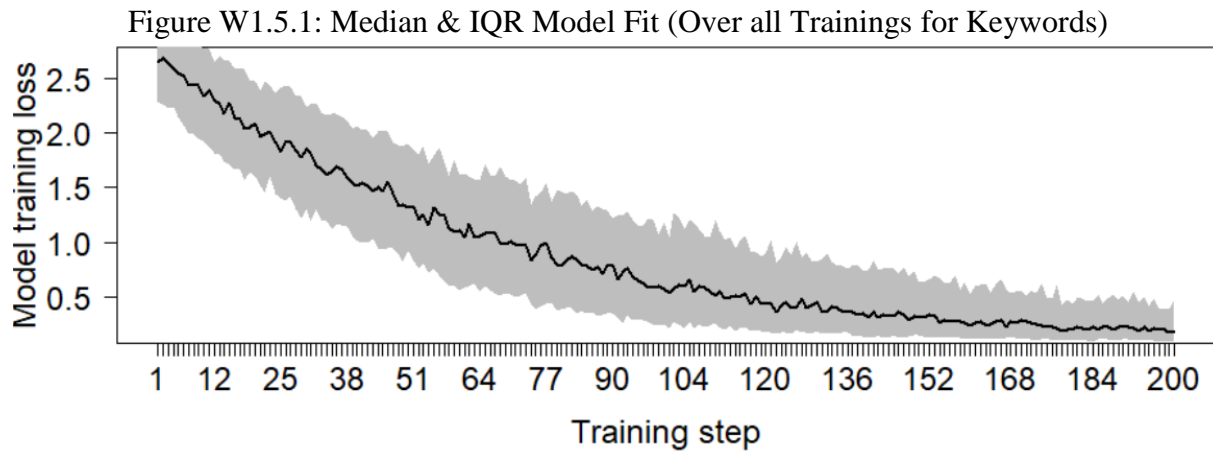
For our experiments, we fine-tune our model for 200 training steps for each keyword, generating 100 pieces of content at each 20th training step which resulted in 1,000 generated texts per focal keyword, of which our method then selected the best scoring pieces of content using the proposed quality score metric. Similar to the approach taken by Liu and Toubia (2018), based on

prior literature and on several test runs, we set the hyper-parameters $\text{top_}k = 40$, and temperature = 0.7 (which effectively regulates the randomness in GPT-2’s sampling process and output content). Next, we show that fine-tuning for 200 training steps is sufficient and examine factors that determine at which training step our proposed method selects the most optimal content.

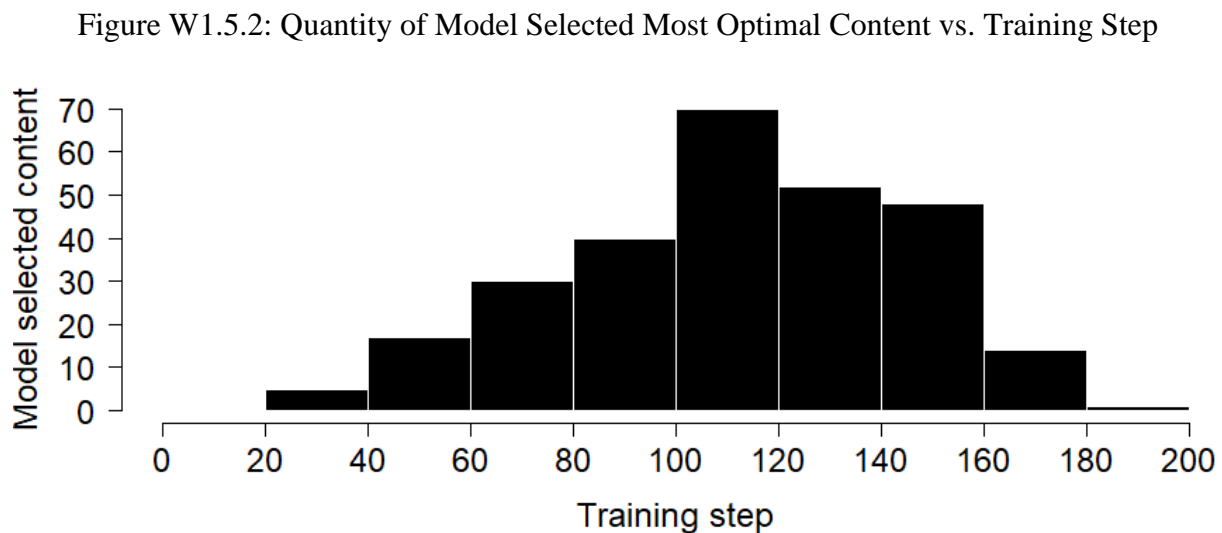
Figure W1.5.1 illustrates the increasing capability of the model to accurately predict words given prior word sequences over the 200 model training steps using the median (black line) and IQR (grey area) of the Loss measure (Radford et al. 2018) over all keyword trainings for the experiments presented in the manuscript. While model fit is consistently improving, Figure W1.5.2 shows that the most optimal content on the basis of the quality score commonly comes from mid training steps (between 60 and 160), while an extremely low and an extremely high amount of training steps entail a lower probability to produce the optimal content. Thus, using 200 training steps for fine-tuning is sufficient.

Using a robust regression (robust against violations of classic data assumptions of regression, see Maechler et al. 2020) for the training steps generating the “best” texts with highest overall quality scores on the quality score components, we observe in Table W1.5 that the content uniqueness among the top 10 ranked websites is the most important determinant for at which training step the most optimal content is generated. That means that when the top 10 ranked websites are more unique compared to each other (i.e., the top ranked websites on which we fine-tune do not make use of many common phrases) our method selects content from a later training phase ($B=117.88$, $t=4.47$, $p<.000$). This may arise because the risk to pick up the repetitive language patterns is lower, and additional fine-tuning steps are needed because the top search results contain more unique phrases. Interestingly, the regression model explains just ~11% of the variance in the data ($\text{Adj.R}^2=.1084$), meaning that the probabilistic fine-tuning and text generation

processes of the GPT-2 model has a considerable impact on at which training step the most optimal content is generated.



— Median of model training loss for all model trainings; ■ IQR of model training loss for all model trainings



■ Quantity of mean top model selected content (for each keyword, we extracted the top scoring generated content and calculated the mean training step from which these came from)

Table W1.5: Quality Score Factors Determining the Training Step for Optimal Content Selection

Robust Regression¹				
Independent Variables	<i>B</i>	Std. Error	<i>t</i>	<i>p</i>
Intercept	54.79	26.56	2.06	.039*
Topic (s_a) + Keywords (s_k) of Top 10	4.82	6.92	0.69	.486
Uniqueness (s_d) of Top 10	117.88	26.37	4.47	<.000**
Readability similarity (s_r) + Naturality similarity (s_n) of Top 10	-31.07	9.59	-3.24	.001**
Adjusted R^2 of regression model: .1084				

¹Dependent variable: Model training step at which most optimal content was selected based on quality score; statistical significance codes: *0.05 level, **0.01 level; because of strong pairwise correlations, we combined s_a and s_k as well as s_r and s_n into one variable by adding them up.

1.6 External Validation of Method Performance

In this section, we assess the generalizability of our proposed method across keywords and industries using our quality score measure. For this purpose, we randomly choose 338 keywords from the approximately 8,500 keywords used previously (typically 9 or 10 keywords for each of the 36 industries) and generated 338,000 pieces of content (1,000 for each single keyword), of which the method automatically selected the best scoring 338 texts (1 for each keyword). Descriptives are in Table W1.4.1, columns (3) and (4).

Table W1.6 reports the difference in medians between the machine generated content and the top 10 ranked websites for all five quality score components in bold, with Wilcoxon rank sum group comparison tests as a statistical difference indicator. We find that the raw machine outperforms the top 10 ranked content for most quality score components in all four industry sectors (Table W1.6). For example, our method outperforms the top 10 ranked websites in terms

of topic consistency (s_a) by ~9% in industry sector I (+.09**), scoring at 34% in topic consistency. The uniqueness of the generated content (s_d), is the only quality indicator that shows a slightly lower value in comparison to the top 10 ranked websites (e.g., -.03** (-3%) in industry sector III), though being at a high value in absolute terms (e.g., ~87% in industry sector III).

Table W1.6: Machine vs. Top 10 Quality Score (All Industry Sectors)

Industry Sector	Statistics	Topic (s_a) ¹	Keywords (s_k) ¹	Uniqueness (s_d) ¹	Readability similarity (s_r) ¹	Naturality similarity (s_n) ¹
I.	Raw Machine vs. Top 10¹ Raw Machine Median ²	+.09** .34	+.14** .34	+.03* .88	+.31** .91	+.25** .83
II.	Raw Machine vs. Top 10¹ Raw Machine Median ²	+.08** .40	+.13** .40	-.02 .88	+.22** .83	+.24* .83
III.	Raw Machine vs. Top 10¹ Raw Machine Median ²	+.10** .43	+.14** .44	-.03** .87	+.22** .83	+.07 .67
IV.	Raw Machine vs. Top 10¹ Raw Machine Median ²	+.11** .40	+.15** .40	-.04* .88	+.31** .91	+.23** .83

¹ Difference in quality score component median value: raw machine generated content vs. real top 10 ranked websites; p-value from Wilcoxon rank sum 2-group comparison tests between machine generated content and top 10 ranked websites; statistical significance codes (one-tailed): *0.05 level, **0.01 level;

² Median quality score component value for raw machine generated content; n=338;

1.7 Alternative Quality Score Weighting Performance Assessment

Our quality score (qs_g) currently consists of 5 dimensions that are weighted equally as depicted in formula (1) in the main manuscript (i.e., the 5 dimensions are just multiplied to get

qs_g). To evaluate the current weighting as depicted in formula (1) in the main manuscript, we consider an alternative weighting scheme: maximize the products of s_a , s_k , s_d , while having minimum thresholds (i.e., cutoff-values) for s_n and s_r . Table W1.7 compares the scores of the current against the alternative weighting scheme for all content produced for our field experiments in the IT service industry and the education sector when applying a 50% (i.e., keep 50% of the top scoring pieces of content) and a 25% (i.e., keep 25% of the top scoring pieces of content) cutoff value for s_n and s_r . A positive (negative) value in Table W1.7 means the quality score weighting scheme presented in the article according to formula (1) performs better (worse) than the alternative quality score weighting scheme. For example, using a cutoff value of 50% in the IT service industry sector experiment, the proposed quality score weighting scheme is superior for s_a (.022**), s_k (.042**), s_d (.042**), and qs_g (.008**), which is consistent across both the used cutoff values and experimental contexts. Thus, the alternative quality score weighting we considered does not result in any improvement, as cutting off content with lower s_n and s_r often results in discarding content that performs well in terms of s_a , s_k , s_d , which ultimately translates into a lower score for the overall quality metric qs_g .

Table W1.7: Comparison of Current vs. Alternative Quality Score Weighting Scheme

Field experiment	s_n & s_r cut-off value ¹	s_a^2	s_k^2	s_d^2	s_n^2	s_r^2	qs_g^2
IT service	50%	.022**	.042**	.042**	-.083**	-.021**	.008**
IT service	25%	.022**	.046**	.178**	-.167**	-.043**	.022**
Education	50%	.030**	.045**	.013**	-.083**	-.043**	.003**
Education	25%	.045**	.059**	.021**	-.083**	-.021**	.014**

¹The cut-off value specifies how many top-scoring data points to maintain, i.e., 50% means keep 50% of the top scoring data-points in s_n & s_r , 25% means keep 25% of the top scoring data points in s_n & s_r (25% is thus more conservative)

²Reported values are median difference values, i.e., median qs_g score of old quality score scheme minus median qs_g score of new (as suggested by the reviewer) quality score scheme. A positive value means the old scheme is superior, a negative value means the new scheme is superior. Significance codes come from two-tailed Wilcoxon rank sum 2-group comparison tests: *0.05 level, **0.01 level;

1.8 Machine Generated Sample Content

To demonstrate the versatility of our approach, Tables W1.8.1, W1.8.2, and W1.8.3 present abbreviated samples of machine-generated content for keywords from varied industries (“best e bike insurance,” “aerospace component manufacturer,” and “state library bookshop”) that have not yet been revised by a human. For comparison reasons each of Tables W1.8.1, W1.8.2, and W1.8.3 contain a real example of top 10 and worst 10 search engine ranked content (i.e., in the examples presented here we used Google’s search engine rankings 290-300). As in Table 2, the texts are shown with their associated quality score components. We can see that while lower ranked real content is both off-topic and performs poor in terms of quality score components (except for uniqueness, for reasons as already discussed in section 1.4 of this

Appendix), the machine generated content typically matches those of the top 10 samples very well and points to subtle differences not immediately transparent to human readers.

Table W1.8.1: Example Generated Piece of Content vs. Top & Worst ranking Content

Stage	Quality Score	Generated Content
GPT-2 fine-tuned (not revised)	$s_a = .65$	Best e bike insurance
	$s_k = .73$	Best value electric bike insurance
	$s_d = .84$	We offer a range of products to suit every budget and circumstance. Our products can provide additional protection for theft, damage, accident, negligence, fraud and more. Our customer service reps are here to tailor a program to your needs.
	$s_n = .58$	
	$s_r = .62$	
	$qs_g = .144$	
		Bicycle Roadside Assistance
		If you find yourself with a disabled bike well come pick you up. [...]
Real top10 ranked content	$s_a = .58$	The importance of e-bike insurance
	$s_k = .61$	Electric bikes (commonly referred to as e-bikes) are expensive, high potential risks on the road typically not fully covered by homeowners, renter or auto insurances. Its a risky move for your e-bike to be uninsured. Fortunately, our bicycle insurance is designed specifically for e-bikes to give you the peace of mind you deserve. Lets take a closer look at why e-bike insurance is a good move. E-bikes can be expensive. [...]
	$s_d = .96$	
	$s_n = .50$	
	$s_r = .79$	
	$qs_g = .133$	
Real worst10 ranked content	$s_a = .19$	Cycle brands in mumbai
	$s_k = .08$	And summer seasons are the perfect time to shape yourself into an envious body, dont you think so? You would be thinking why are we talking [Continue reading] Top 20 Sportswear Brands for Men and Women in India Metro Cycle Company in Charni Road, Mumbai. Know more about latest cycles and new cycles in India 2020. [brand name] [brand name] ([Phone number]) is the leading Bicycle Dealer In India. Incense sticks nowadays are made up of aromatic materials that produce aroma when burned. Cycles in India Get all Cycles information like cycle prices, reviews, photos, etc. . There are more cycle brands in India but we have listed only 25, drop your comment for your best cycle. , one of the largest cycle manufacturers in India is now emerging as best Bicycle Companies online. [...]
	$s_d = .96$	
	$s_n = .17$	
	$s_r = .21$	
	$qs_g = .001$	

Human revision in our field experiment reported below: ■ = identifiers like brand names and phone numbers are replaced by a tag (e.g., [brand name]) to retain confidentiality; headlines are printed in bold to ease reading;

Table W1.8.2: Example Generated Piece of Content vs. Top & Worst ranking Content

Stage	Quality Score	Generated Content
GPT-2 fine-tuned (not revised)	$s_a = .40$	Aerospace component manufacturer
	$s_k = .40$	
	$s_d = .88$	With a revenue of \$8 billion, [brand name] is the third-largest aerospace company in the world. It provides products for both global aerospace and defense, manufacturing nacelle systems, avionics systems, control surfaces, electrical systems, mechanical and aerospace systems, and flight simulator systems.
	$s_n = .92$	
	$s_r = .98$	The company also provides aftermarket services, like draping and electrical systems, in addition to product development.
	$qs_g = .127$	[brand name] specializes in fabricating complex components and sub-assemblies for military and commercial uses. Examples of products it creates are for use in communication systems, missiles and aviation defense, and in the aerospace industry generally. [...]
Real top10 ranked content	$s_a = .31$	[brand name]: World Leaders in Aerospace
	$s_k = .38$	
	$s_d = .99$	Since Flight Began [brand name] comprises one of the largest, most experienced concentrations of world-class aerospace companies, the world's AEROSPACE ALLEY! Our firms grew up where aerospace was born. Today, our precision components take to the skies every day, in every corner of the globe, as they have since the inception of powered flight.
	$s_n = .25$	
	$s_r = .60$	[brand name] member companies are united by a single goal; be the world leader in providing customers with aerospace components of unsurpassed quality, at competitive prices and always, on-time delivery. To achieve that goal, we work with our peers and with regional, national and international experts to ensure continuous improvement of our products and services. [...]
	$qs_g = .017$	
Real worst10 ranked content	$s_a = .07$	Aerospace Industry
	$s_k = .05$	
	$s_d = .91$	Aerospace Industry
	$s_n = .17$	
	$s_r = .13$	TOTAL: 8 HORNET Series INSPIRER Series GRANDER 5Ax Series DBC2000mm DBC1500mm GRANDER 5MG Series DBC3100mm DBC2500mm Industry Aerospace Industry GRANDER 5Ax Series GRANDER 5MG Series HORNET Series INSPIRER Series DBC 1500mm DBC 2000mm DBC 2500mm DBC 3100mm Automobile HE Series HT Series VTW Series VA Series Medical Sphere VTH Series VTP Series VTJ Series VTG Series Large Molding DBC 1500mm DBC 2000mm DBC 2500mm DBC 3100mm Electronics Products HF Series VA Series VF Series VTT Series Parts Processing HE Series HF Series VTT Series VTW Series VF Series VA Series VH Series VE Series VK Series VP Series VTH Series TS Series CNC Lathes T Series CNC Lathes TEL: [phone number] FAX: [fax number] Email: [email] Address: [address] TEL: [phone number] Email: [email] Address: [address] 2018 [brand name] All rights reserved.
	$qs_g = .000$	

Human revision in our field experiment reported below: [] = identifiers like brand names and phone numbers are replaced by a tag (e.g., [brand name]) to retain confidentiality; headlines are printed in bold to ease reading;

Table W1.8.3: Example Generated Piece of Content vs. Top & Worst ranking Content

Stage	Quality Score	Generated Content
GPT-2 fine-tuned (not revised)	$s_a = .25$	State library bookshop
	$s_k = .22$	open seven days a week
	$s_d = .88$	[...]
	$s_n = .92$	\$10 kids 12 and under
	$s_r = .98$	Free playtime on the playroom terrace
	$qS_g = .042$	Discover the wonders of Queensland
		The [brand name] Historical Library is Queensland's first permanent, accessible, fully-staffed library. The Library is home to a beautifully curated selection of Queensland's best contemporary and classic books, and offers free access to more than 220,000 unique readers each day. The Library is Brisbane's only permanent non-profit library, and is funded entirely by the State Government. [...]
Real top10 ranked content	$s_a = .27$	Readings State Library
	$s_k = .25$	Readings State Library is the first permanent commercial bookshop permitted to be established in the building in the Library's 150 year history. It is a unique space that caters to students, travellers, book lovers, CBD workers and library regulars. As part of the State Library Victoria's Vision 2020 redevelopment project, our shop was relocated to a welcoming new space beautifully designed by Australasian design firm Architectus with Scandinavia's Schmidt Hammer Lassen Architects. [...]
	$s_d = .96$	
	$s_n = .42$	
	$s_r = .72$	
	$qS_g = .019$	
Real worst10 ranked content	$s_a = .11$	About the Author
	$s_k = .03$	Tasmanian writer Katherine Johnson is the author of four novels. Paris Savages based on a true story of human zoos will be published in October 2019. Katherine's third novel, Matryoshka or Russian dolls (Ventura Press 2018), is a story of secrets, refuge and loves lost and found. Her previous novels include The Better Son (Ventura Press 2016), set in northern Tasmanias caves, and Pescadors Wake (Fourth Estate 2009), the story of the danger and heartbreak of lives at the mercy of the sea during a three-week Southern Ocean chase. Katherine is the recipient of The University of Tasmania Prize and the Peoples Choice Award (Tasmanian Literary Prizes) and HarperCollins Varuna Awards. The Better Son was longlisted for the Australian Indie Book Awards 2017 and the Tasmania Book Prize (Premiers Literary Prizes 2017). [...]
	$s_d = .96$	
	$s_n = .25$	
	$s_r = .74$	
	$qS_g = .000$	

Human revision in our field experiment reported below: [] = identifiers like brand names and phone numbers are replaced by a tag (e.g., [brand name]) to retain confidentiality; headlines are printed in bold to ease reading;

2 IT Service Industry Field Study

This section contains details and accompanying materials for the empirical performance evaluation of our semi-automated content generation machine in the IT service industry application setting reported in the article.

2.1 Keywords Used & Keyword Statistics

Table W2.1 depicts the keywords (search queries) for which the experimental groups in the IT service sector field study produced content, including keyword statistics and descriptive statistics for the ranking performance of the revised machine content in the search engine. The basic keyword statistics reported in Table W2.1 include the average monthly search volume that serves as an indicator of how many users on average search for the keyword per month, the paid keyword competition to capture SEA competition that is provided by the search engine, and the keyword length to account for how many sub-words constitute the keyword. Our company partner selected the keywords used in the experiment based on its standard procedure for keyword selection (i.e., based on monthly search volume, competition, fit with the firm and keyword strategy), with a preference in favor of keywords with lower search volume in the long tail as part of their keyword strategy. The selected keywords are similar in terms of keyword statistics (competition, search volume, keyword length), and the company did not have any prior search engine ranking history for any of the keywords that they provided to us.

Table W2.1: Keywords for IT Service Field Experiment

Keyword	Descriptives								
	Avg. monthly search volume	Competition	Competition index	Keyword length	Mean revised machine ranking	Median revised machine ranking	SD revised machine ranking	IQR revised machine ranking	% of days revised machine was in ranking
IT procurement	10	low	4	2	11.65	8	7.68	10	90.70
IT support and services	10	low	3	4	15.61	15	6.54	5	97.21
global IT support	10	-	-	3	15.03	14	5.11	3	72.09
IT assessment	10	low	0	2	21.58	22	5.87	6	87.91
IT consulting services	10	-	-	3	18.22	17	7.43	5	66.05
IT maintenance	10	-	-	2	13.66	12	8.11	11	99.07
IT service maintenance	0	-	-	3	3.19	2	3.25	2	99.07
IT service support	10	low	0	3	10.74	10	2.82	4	99.07
IT service continuity	10	-	-	3	21.57	19	6.47	9	92.56
IT support business	10	-	-	3	51.60	52	14.29	26	39.53
Small business IT support services	0	-	-	5	94.75	50	79.96	129	64.19
IT support costs for small business	0	-	-	6	14.21	13	3.96	6	99.07
IT maintenance support	0	-	-	3	4.11	2	3.86	6	99.07
IT maturity assessment	10	low	29	3	22.67	22.5	7.10	9.75	97.67
IT procurement services	10	-	-	3	3.23	3	1.11	2	99.07
IT procurement process	10	-	-	3	30.18	19	22.22	21	99.07
IT solution delivery	10	low	0	3	2.47	1	2.21	2	90.70
IT strategy consulting	10	-	-	3	28.25	26	9.24	5	98.14
IT consulting software	10	high	100	3	15.70	15	5.16	8	75.35

Entries that display “-“ mean that the search engine keyword tool did not provide specific information.

2.2 Experimental Setup

To generate website landing page content for the company we collaborated with, the participants of the three groups of human content writers described in the IT service industry application study in the article had access to all tools and the environment that the company uses

in its common SEO content production workflow. The instructions provided to the study participants are presented in Table W2.2.1. The participants were offered incentives for content production. The incentive for groups 1 (novices) and 2 (quasi-experts) was 15 € per produced content and credit for a marketing course. The incentive for group 3 (SEO experts) was 40 € per produced content.

Content production took place within the same week and in the same geographic location so that all participants had the same state of search engine results as a basis, which we controlled for via daily crawls. We controlled for writers' differences in terms of educational and other background covariates we considered relevant for the writing task. To this end we conducted a series of Kruskal Wallis group comparison tests. We find that the human content writing groups did not differ in their education ($\chi^2(3)=.60$, $\eta^2=.01$, $p=.745$) or writing skills (for which the SEO experts scored a bit higher; $\chi^2(3)=5.89$, $\eta^2=.12$, $p=.053$); we also could not detect any significant differences in terms of time invested conducting research on the target keyword / topic ($\chi^2(3)=.28$, $\eta^2=.00$, $p=.868$) and content writing ($\chi^2(3)=3.76$, $\eta^2=.08$, $p=.153$).

Table W2.2.1: Participants' Survey Instructions for Content Writing

Content Writing Group	Instructions ¹
Novices	<p>[Short introduction stating the goal of this study, strict anonymization, the incentive and a contact person for questions.]</p> <p>Imagine you are a marketing employee in an IT service company.</p> <p>Your manager approaches you to write a Google search engine optimized (SEO) text for a single site on the website of your IT company, that elaborates on a specific service. You should write the text in a way that it ranks well in Google. That means, it should preferably appear on page 1 in the Google search results.</p>

- The text should be written for the keyword / search term / topic: “**IT maintenance**” (i.e., for IT maintenance provided as a service by your company to firms).
- It should be written for ranking well in **Google in [Country blinded], set to English language** (please use the link below).
- For ranked **example sites see:**
<https://www.google.com/search?num=100&hl=en&q=it+maintenance>
- It should be **original, unique content**, invented by you (i.e., NO copies).
- It should be written **in English language**.
- It should contain **around 700 to 800 words** (ca. 2 A4 pages).

Your text: (Please write your text in the following text field.)

Quasi Experts

[Short introduction stating the goal of this study, strict anonymization, the incentive and a contact person for questions.]

Imagine you are a marketing employee in an IT service company.

Your manager approaches you to **write a Google search engine optimized (SEO) text for a single site on the website of your IT company**, that elaborates on a specific service. You should **write the text in a way that it ranks well in Google**. That means, it should preferably appear on page 1 in the Google search results.

- The text should be written for the keyword / search term / topic: “**IT maintenance**” (i.e., for IT maintenance provided as a service by your company to firms).
- It should be written for ranking well in **Google in [Country blinded], set to English language** (please use the link below).
- For ranked **example sites see:**
<https://www.google.com/search?num=100&hl=en&q=it+maintenance>
- It should be **original, unique content**, invented by you (i.e., NO copies).
- It should be written **in English language**.
- It should contain **around 700 to 800 words** (ca. 2 A4 pages).

How to write a SEO optimized text?

- **Integrate the main keyword** (“IT maintenance”) **or parts of it most often compared to the other words in your text.**
- **Write about subtopics / content** that you can find on the top ranked websites for the main keyword.

- **Align the word distribution of your text** with the word distribution of the top ranked websites for the main keyword (i.e., **put the right words with the right frequencies into your text**).
- For the **word distribution analyses** use: <https://wordcounter.net/> (Please be aware that the tool doesn't count common stopwords like "it".)
- **Prevent keyword stuffing** (i.e., [don't integrate keywords overly often and in an unnatural way into your text](#)).
- Try to give your text a **good readability and structure**.

Your text: (Please write your text in the following text field.)

Real SEO
Experts

[Short introduction stating the goal of this study, strict anonymization, the incentive and a contact person for questions.]

Imagine you are a marketing employee in an IT service company.

Your manager approaches you to **write a Google search engine optimized (SEO) text for a single site on the website of your IT company**, that elaborates on a specific service. You should **write the text in a way that it ranks well in Google**. That means, it should preferably appear on page 1 in the Google search results.

- The text should be written for the keyword / search term / topic: **“IT maintenance”** (i.e., for IT maintenance provided as a service by your company to firms).
- It should be written for ranking well in **Google in [Country blinded], set to English language** (please use the link below).
- For ranked **example sites** see: <https://www.google.com/search?num=100&hl=en&q=it+maintenance>
- It should be **original, unique content**, invented by you (i.e., NO copies).
- It should be written **in English language**.
- It should contain **around 700 to 800 words** (ca. 2 A4 pages).

Your text: (Please write your text in the following text field.)

¹Keywords and links were adapted in each survey.

Descriptive statistics on the content length and changes are provided in Table W2.2.2. As these statistics show, the produced content is of about equal length across the experimental

groups and human revisers changed about 9% of the machine-made content in the revision process. We report evaluations of the performance of each of the experimental groups with regards to search engine rankings and the quality score in the next section of this document.

Table W2.2.2: Descriptives for Content Lengths and Revision

Dimension	Groups	Descriptives			
		Median	(IQR)	Min	Max
Produced content length (in words)	Revised machine	807	(67)	632	899
	Real SEO Experts	729	(84)	578	771
	Quasi Experts	694	(69.5)	498	749
	Novices	711	(48.5)	377	966
Content change (raw vs. revised)¹	Change in %	9.04	(3.77)	3.31	21.45
	Change in words	74.00	(36.50)	27.00	154.00

¹This includes every possible change between the raw machine and revised machine output like added words, deleted words, and words with at least one changed letter (including changed letter capitalization).

2.3 Performance of the Experimental Groups

2.3.1 Quality Scores of Experimental Groups' Content

In Table W2.3.1, we compare the quality score components for the content generated by each of the experimental groups and the top 10 ranking search results. The topic (s_a), keyword (s_k), and readability similarity (s_r) scores are higher for the raw and semi-automated content compared to the remaining experimental groups and the top 10 ranked websites as well as the

lowest ranked search results, while human created content scores higher in terms of uniqueness. This finding demonstrates the content writing machine’s capability to mimic patterns such as topics, keywords and readability levels contained in the top10 pages more thoroughly than human writers.

In comparing the experimental groups to the top 10 search engine rankings, bear in mind that the scores are comparative measures. For most of the components s_a , s_k , s_r and s_n , the top 10 ranked pages score lower than the raw and revised machine-generated content (and also some of the groups of human writers). This means that the machine-generated content is more similar along these dimensions to the set of the top ranked pages than the group of top 10 ranked pages are internally (i.e., comparing any given page from the top ranked pages to the set of top ranked pages). For example, consider the search query / keyword “IT service management”. Content ranked on places 1-5 might focus about a relevant sub-topic (e.g., a telephone hotline for IT service management) of the keyword / search query, and content ranked on places 6-10 might focus on a different relevant sub-topic (e.g., a ticketing system for better customer support in IT service management) for the keyword / search query. The content writing machine is opportunistic in a sense, as it taps into features from all of the top 10 pages. That is, it picks, combines and rearranges features from all of the top 10 pages in a novel way. This “averaging” effect occurs during the fine-tuning process, which makes use of all of the top 10 pages. This “averaging” is taken into account in the quality score, in which all components compare a given content to the set of the top ranked content, which is the essence of SEO: trying to integrate all relevant features and aspects (that are present in the content of the top ranked pages) to fulfill the user’s search query best (e.g., to provide the most rich information, formulated in the most appropriate way) (see section 1.5 in this document).

Table W2.3.1: Quality Score Components Group Comparisons to Top 10 Ranked Websites

Quality Score Component	Group	Descriptives				Wilcoxon rank sum ¹			
		Median (IQR)	Min	Max		W	z	r	p
Topic (<i>s_a</i>)	Top 10	.32	(.11)	.11	.27				
	Revised machine	.40	(.13)	.35	.68	65	3.36	.54	.000**
	Raw Machine	.46	(.13)	.33	.61	62	3.60	.58	.000**
	Real SEO Experts	.37	(.08)	.30	.49	59	1.27	.24	.205
	Ouasi Experts	.36	(.10)	.10	.61	139	1.19	.19	.234
	Novices	.29	(.08)	.20	.56	211	-.87	-.14	.385
	Worst 10	.19	(.07)	.11	.28	344	-5.39	-.87	.000**
Keywords (<i>s_k</i>)	Top 10	.30	(.16)	.05	.26				
	Revised machine	.44	(.18)	.32	.74	68	3.27	.53	.001**
	Raw Machine	.48	(.18)	.31	.62	64	3.53	.57	.000**
	Real SEO Experts	.36	(.06)	.16	.51	67	.87	.16	.383
	Ouasi Experts	.38	(.13)	.01	.70	142	1.10	.17	.271
	Novices	.31	(.22)	.52	.61	213	-.93	-.15	.354
	Worst 10	.16	(.11)	.05	.26	335	-4.97	-.81	.000**
Uniqueness (<i>s_d</i>)	Top 10	.92	(.04)	.79	.97				
	Revised Machine	.90	(.06)	.81	1.00	216	-1.02	-.16	.307
	Raw Machine	.84	(.12)	.52	.94	301	-3.66	-.59	.000**
	Real SEO Experts	.98	(.03)	.93	1.00	15	3.45	.65	.000**
	Ouasi Experts	.99	(.03)	.86	1.00	58	3.56	.58	.000**
	Novices	.98	(.07)	.79	1.00	86	2.74	.45	.006**
	Worst 10	.95	(.04)	.89	.98	108	2.11	.34	.034*
Readability Similarity (<i>s_r</i>)	Top 10	.56	(.08)	.50	.74				
	Revised Machine	.87	(.17)	.47	1.00	21	4.64	.75	.000**
	Raw Machine	.96	(.09)	.70	1.00	2	5.21	.84	.000**
	Real SEO Experts	.57	(.51)	.21	1.00	77	.39	.07	.694
	Ouasi Experts	.53	(.39)	.10	1.00	206	-.73	-.12	.465
	Novices	.57	(.42)	.08	.96	176	-.11	-.02	.907
	Worst 10	.47	(.10)	.26	.68	312.5	-3.84	-.62	.000**
Naturalness Similarity (<i>s_n</i>)	Top 10	.56	(.05)	.50	.63				
	Revised Machine	.67	(.38)	.17	1.00	145	1.02	.16	.306
	Raw Machine	.92	(.38)	.42	1.00	55.5	3.66	.59	.000**
	Real SEO Experts	.75	(.33)	.33	1.00	25	2.96	.56	.003**
	Ouasi Experts	.75	(.38)	.17	.83	132.5	1.39	.23	.164
	Novices	.58	(.38)	.00	1.00	162	.53	.09	.598
	Worst 10	.35	(.15)	.18	.53	359	-5.20	-.84	.000**

¹Two-tailed tests; statistical significance codes: *0.05 level, **0.01 level;

2.3.2 Post Hoc Test for Achieved Search Engine Rankings

In Table 2 of the article we reported that the semi-automated content outperforms the remaining experimental groups in terms of (top 10) search engine rankings. To check robustness of this finding and to test whether it also generalizes for pairwise comparisons of experimental groups we conducted a series of Kruskal Nemenyi post hoc tests. Table W2.3.2 reports the results which show that the search engine performances of all experimental groups are statistically different at the 0.05 level.

Table W2.3.2: Post Hoc Tests: Search Engine Rankings Performance Comparison (IT Service Sector)

Dimension	Group	Kruskal Nemenyi Post Hoc Test (<i>p</i>)		
		Real SEO Experts	Quasi Experts	Novices
Pages in ranking / day	Revised Machine	<.000**	<.000**	<.000**
	Real SEO Experts		.014*	<.000**
	Quasi Experts			<.000**
Pages in top 10 / day	Revised Machine	<.000**	<.000**	<.000**
	Real SEO Experts		.003**	<.000**
	Quasi Experts			<.000**
Mean rankings / day	Revised Machine	<.000**	<.000**	<.000**
	Real SEO Experts		<.000**	<.000**
	Quasi Experts			<.000**

¹Statistical significance codes: *0.05 level, **0.01 level, chi-square approximated;

2.3.3 Content Ranking for Sub-Keywords Assessment

In SEO, content is usually optimized for a single main keyword (search query).

However, in search engine advertising (SEA) ads and bids are often optimized for multiple keywords at once. Thus, a company could potentially benefit from optimizing SEO content for multiple keywords simultaneously. That is why we conducted an additional experiment for our IT service industry study to explore how well the experimental groups' SEO content ranks for related keywords. We identified the latter by analyzing the word distributions of the top 10 search engine ranked content for the 19 main keywords specified in Table W2.1 and extracted the most frequent keywords and groups of words, yielding 207 related keywords. For example, when analyzing the top 10 search results for the keyword "IT assessment", we find the (most frequently occurring) following related keywords based on their word distributions: "IT assessments", "business continuity", "disaster recovery", "security assessment", "assessment services", "IT assessment services", "risk security assessment", "information technology assessment", "disaster recovery plan". After scraping the search engine rankings for these 207 related keywords, we find that the revised content machine ranks substantially better and more often for related keywords than the competing human groups (Table W2.3.3). For example, the revised machine ranked for 34 related keywords and occurred in the top 10 results six times. The median search engine ranking of the revised machine is 23.

Based on this auxiliary study, we conclude that the method seems to perform surprisingly well for related keywords as well. In contrast to relying on heuristics such as keyword density, the fine-tuning process of our semi-automated algorithm appears to not only capture the overall topic but also related sub-topics within the content. Thus, in the process of

generating content for a specific keyword, our content also performs reasonably well in terms of search engine rankings for topic-related sub-keywords for which it was not primarily optimized.

Table W2.3.3: Ranking Performance of Content for Related Keywords

Group	Descriptives			
	Median (IQR) search engine ranking		Total number of ranked pages	Number of pages ranked in top10
Revised machine	23	(23.00)	34	6
Real SEO Experts	26	(35.50)	4	1
Quasi Experts	188.5	(78.25)	4	0
Novices	68	(27.50)	3	0

Desriptives for achieved rankings per experimental group for topic-related sub-keywords extracted from the top 10 ranked pages (207 sub-keywords, and 51,995 total ranked pages).

2.3.4 Additional Keywords Performance

In this section, we assess the content machine’s performance for additional keywords in the IT service sector experiment, which were not included in our original study, bringing the keyword count for all experimental groups to 30 keywords. Due to technical issues with the company’s website that were beyond our control, we were unable to put the generated content online to evaluate search engine rankings. Nonetheless, we report the quality scores for the machine and human made content.

Table W2.3.4 depicts the additional keywords for which content was generated. Consistent with our previously reported findings, Table W2.3.5 illustrates that the raw and revised machine content substantially outperforms all competing human content producing groups including the SEO experts, similar to the results depicted in Table W2.3.1.

Table W2.3.4: Additional Keywords for the IT Service Field Experiment

Field Study	Keyword	Descriptives			
		Avg. monthly search volume	Competition	Competition index	Keyword length
IT service	SLA contract	10	low	0	2
	SLA ITIL	10	low	14	2
	service level agreement best practices	10	low	0	5
	IT security services	10	-	-	3
	ITIL incident	20	low	11	2
	ITIL ITSM	10	low	32	2
	IT maintenance contract	10	-	-	3
	IT project management	40	low	18	3
	IT scalability	10	low	0	2
	IT performance management	10	low	26	3
	Server maintenance	20	low	0	2

Entries that display “-” mean that the search engine keyword tool did not provide specific information.

Table W2.3.5: Quality Score Components Group Comparisons to Top 10 Ranked Websites (Keyword Count Increased to 30 Keywords)

Quality Score Component	Group	Descriptives				Wilcoxon rank sum ¹			
		Median (IQR)		Min	Max	W	z	r	p
Topic (<i>s_a</i>)	Top 10	.38	(.23)	.21	.69				
	Revised machine	.49	(.22)	.35	.71	253	2.91	.38	.004**
	Raw Machine	.50	(.19)	.28	.77	260	2.84	.37	.004**
	Real SEO Experts	.40	(.19)	.30	.72	302	2.01	.26	.043*
	Ouasi Experts	.41	(.19)	.10	.68	381	.81	.10	.420
	Novices	.33	(.16)	.17	.64	490	-1.08	-.14	.281
	Worst 10	.18	(.07)	.09	.28	870	-7.23	-.93	.000**
Keywords (<i>s_k</i>)	Top 10	.38	(.26)	.14	.77				
	Revised machine	.52	(.28)	.32	.81	250	2.95	.38	.003**
	Raw Machine	.55	(.23)	.29	.85	260	2.84	.37	.004**
	Real SEO Experts	.46	(.21)	.16	.79	318	1.77	.23	.077
	Ouasi Experts	.43	(.21)	.01	.77	382	.79	.10	.429
	Novices	.33	(.25)	.05	.72	498	-1.20	-.16	.230
	Worst 10	.15	(.10)	.03	.26	853	-6.77	-.87	.000**
Uniqueness (<i>s_d</i>)	Top 10	.92	(.09)	.72	.99				
	Revised Machine	.91	(.13)	.74	1.00	418	.47	.09	.641
	Raw Machine	.81	(.12)	.52	.97	694	-3.70	-.48	.000**
	Real SEO Experts	.97	(.06)	.68	1.00	203	3.51	.46	.000**
	Ouasi Experts	.99	(.03)	.87	1.00	101	5.07	.66	.000**
	Novices	.99	(.04)	.79	1.00	161	4.02	.53	.000**
	Worst 10	.94	(.05)	.85	.99	323	1.87	.24	.061
Readability similarity (<i>s_r</i>)	Top 10	.56	(.08)	.48	.74				
	Revised Machine	.82	(.19)	.26	1.00	104	5.11	.66	.000**
	Raw Machine	.95	(.14)	.70	1.00	6	6.57	.85	.000**
	Real SEO Experts	.70	(.53)	.02	1.00	343	1.39	.18	.165
	Ouasi Experts	.49	(.43)	.04	1.00	520	-1.28	-.17	.200
	Novices	.53	(.45)	.02	.96	497	-1.19	-.16	.234
	Worst 10	.43	(.17)	.12	.68	827.5	-5.57	-.72	.000**
Naturalness similarity (<i>s_n</i>)	Top 10	.56	(.06)	.50	.65				
	Revised Machine	.67	(.40)	.17	1.00	384	.97	.13	.332
	Raw Machine	.75	(.42)	.33	1.00	227.5	3.29	.43	.000**
	Real SEO Experts	.67	(.33)	.08	1.00	346.5	1.34	.17	.181
	Ouasi Experts	.50	(.33)	.17	.83	448.5	-.20	-.03	.843
	Novices	.54	(.44)	.00	1.00	453.5	-.51	-.07	.607
	Worst 10	.30	(.15)	.08	.53	896	-6.59	-.85	.000**

¹Two-tailed tests; statistical significance codes: *0.05 level, **0.01 level;

2.4 Providing Quality Score Feedback to Revise Content

To explore the incremental value of using NLG relative to SEO experts with access to the quality score further, we conducted the following additional study. Following the content that was initially produced by the real SEO experts we offered them the opportunity to improve their content by providing them with the quality scores of their initial content and the real search engine ranked top 10 content. The experimental setup mirrors A/B testing and is as follows: The tests were conducted using an online survey (including personal explanations of the task and a Q&A section). The survey contained the task description, the principal investigator's contact details, and an incentive of 40€ for the revision / feedback round per piece of content. The participants were introduced to the quality score and provided with an explanation of each quality score component and how to interpret it. They were provided with the quality score (both on each component and overall) for the content they produced initially. For comparison purposes, they were also provided with the top 10 ranked content for the specific keyword and their associated quality scores. Study participants entered their revised text in an open text field. We extended the original study reported in the main manuscript using 30 keywords for the real SEO experts (instead of just 9), so the testing was conducted for 30 pieces of SEO expert produced content.

Table W2.4.1 shows that the SEO experts changed their original content by 10.24% (~77.50 words), ranging between 12 words changes and 176 word changes. Table W2.4.2 compares the achieved quality scores of the original SEO experts' content to the revised SEO experts' content for each quality score component. We find no statistically significant differences between them, suggesting that the SEO experts were not able to improve the quality of their content, likely due to the associated complexity (i.e., dozens of word distributions, numbers, and

abstract concepts). This suggests that the semi-automated procedure not only reduces the time/cost associated with content production, but also performs better than human experts on tasks involving the generation of content for a specific purpose.

Table W2.4.1: Descriptives for Real SEO Experts Content Revision

Dimension	Groups	Descriptives			
		Median	(IQR)	Min	Max
Produced content length (in words)	Original SEO Expert Content	729.5	(47.25)	587	819
	Revised SEO Expert Content	760.5	(58.75)	546	930
Content change (original vs. revised)¹	Change in %	10.24	(5.84)	1.62	24.24
	Change in words	77.50	(49.25)	12	176

¹This includes every possible change between the original SEO experts and revised SEO experts content such as added words, deleted words, and words with at least one changed letter (including changed letter capitalization).

Table W2.4.2: Quality Score: Original Real SEO Experts Content vs. Revised Real SEO Experts Content

Quality Score Component	Group	Descriptives				Wilcoxon rank sum ¹			
		Median (IQR)		Min	Max	W	z	r	p
Topic (<i>s_a</i>)	Original SEO Experts	.41	(.18)	.30	.72				
	Revised SEO Experts	.44	(.18)	.30	.72	423	.39	.05	.697
Keywords (<i>s_k</i>)	Original SEO Experts	.46	(.20)	.16	.79				
	Revised SEO Experts	.46	(.15)	.16	.80	423	.39	.05	.697
Uniqueness (<i>s_a</i>)	Original SEO Experts	.97	(.06)	.08	1.00				
	Revised SEO Experts	.96	(.07)	.63	.99	560	-1.62	-.21	.105
Readability similarity (<i>s_r</i>)	Original SEO Experts	.72	(.54)	.02	1.00				
	Revised SEO Experts	.66	(.71)	.02	1.00	469	-.27	-.04	.784
Naturalness similarity (<i>s_n</i>)	Original SEO Experts	.67	(.39)	.08	1.00				
	Revised SEO Experts	.58	(.25)	.17	1.00	477.5	-.40	-.05	.687

¹Two-tailed tests between original vs revised real SEO experts quality scores, statistical significance codes: *.05 level, **.01 level;

2.5 Consumers' Content Perceptions

We provide details for the MTurk study presented in the article in which we examine differences in consumer perceptions between the semi-automated and human content. The instructions with which survey participants were presented are shown in Table W2.5.1.

Table W2.5.1: Participants' Survey Instructions

Survey Instructions
Dear study participant
Thank you for participating in our study on SEO & text writing. Your input is vital for us. In the following, besides answering some demographic questions, we will ask you to read and assess 1 text.
It will take you 5 minutes at most to finish the survey.
Please read all questions and the text mindfully and completely , and answer all questions as honestly and spontaneously as possible . Follow your intuition, there are no right or wrong answers .
All information that you provide to us will be strictly treated as anonymous . Thank you for your kind support.
Sincerely, [...]
[New survey page]
Imagine, you are looking for an IT service for your company, and you come across a website with the text below. Please take a look at it.
[Randomized piece of content]
[Questions to assess content]

To assure data quality in our survey-based content perception experiment, we implemented honeypots (for antispam), attention and honesty checks (i.e., reverse coded items

and same questions worded a bit differently), and excluded all surveys with a completion time lower than 1.50 minutes, leaving us with 551 surveys for our analyses. We performed scale reliability checks using Cronbach's Alpha including deleting offset items. Using a series of Kruskal Wallis tests, we assured that participants' characteristics did not differ substantially between the experimental conditions in terms of the time to finish the survey ($\chi^2(3)=3.38$, $\eta^2=.01$, $p=.337$), the participants' gender ($\chi^2(3)=2.00$, $\eta^2=.00$, $p=.572$), the highest completed level of education ($\chi^2(3)=3.08$, $\eta^2=.01$, $p=.380$), age ($\chi^2(3)=.25$, $\eta^2=.00$, $p=.969$), and English reading proficiency ($\chi^2(3)=.41$, $\eta^2=.00$, $p=.939$).

Table W2.5.2 reports operationalizations, references, and scale reliability metrics for the content perception study we conducted.

Table W2.5.2: Operationalizations & Measures of Main Variables for Survey

Variable	Items	Source	Scale Reliability ¹
Readability	Bipolar 5-point scale with following items: “Please indicate whether you perceive the text above as ... <ul style="list-style-type: none"> ● poorly written – well written ● poorly readable – well readable ● not fitting together well – fitting together well ● not understandable – understandable ● not interesting – interesting” 	Pitler and Nenkova 2008	.91
Understandability	Bipolar 5-point scale with following items: “Please indicate whether you perceive the text above as ... <ul style="list-style-type: none"> ● complicated – simple ● unclear – clear ● chaotic – orderly ● illogically arranged – logically arranged ● wordy – concise ● difficult – easy“ 	Kamoen et al. 2013	.88
Credibility	Bipolar 5-point scale with following items: “Please indicate whether the text above is ... <ul style="list-style-type: none"> ● unbelievable – believable ● inaccurate – accurate ● not trustworthy – trustworthy ● biased – not biased ● incomplete – complete” 	Roberts 2010, Flanigan and Metzger 2000	.87
Attitude toward the content	Bipolar 5-point scale with following items: “Please indicate whether you feel that the text above is ... <ul style="list-style-type: none"> ● distant – appealing ● reluctant – inviting ● boring – fascinating ● impersonal – personal ● monotonous – varied ● interesting – uninteresting” 	Kamoen et al. 2013	.89

¹Cronbach’s Alpha with optimized number of items

In addition to the scale items included in Table W2.5.2, we measure perceived content naturality using two items. On bipolar five-point scales, we ask respondents to indicate whether they believe that the content feels artificial vs. feels natural, and machine-made vs. human-made. We also ask two questions to assess future intent. To gauge willingness to further inform, we use a slider from 0 to 100 and ask respondents to indicate how they agree with the statement: “I want to further inform myself about the company providing the service.” To measure willingness to buy, we use a slider from 0 to 100 and ask respondents to indicate how much they agree with the statement: “I am willing to buy the described service.”

In Table W2.5.3, we report pairwise correlations between user perception variables using Kendall’s tau b, illustrating high correlations between these items.

Table W2.5.3: Consumer Content Perception: Dimensions’ Inter-Correlations

Dimension	Kendall’s tau b (τ_b)						
	Readability	Understandability	Credibility	Attitude Toward the Content	Content Naturality	Willingness to Further Inform	Willingness to Buy
Readability	1.00**	.59**	.57**	.50**	.52**	.41**	.42**
Understandability		1.00**	.43**	.58**	.57**	.44**	.46**
Credibility			1.00**	.40**	.44**	.33**	.37**
Attitude Toward the Content				1.00**	.58**	.52**	.53**
Content Naturality					1.00**	.44**	.49**
Willingness to Further Inform						1.00**	.69**
Willingness to Buy							1.00**

¹Statistical significance codes: *0.05 level, **0.01 level, one-tailed; n=551;

We complement the MTurk Study with a computational linguistic analysis. Using LIWC (Pennebaker et al. 2015), the evaluative lexicon (Rocklage et al. 2018), and the text analyzer

(Berger et al. 2020b) software packages that apply various lexica, analyses and scales, we assess the linguistic properties along psychological dimensions including concreteness, familiarity, and emotionality. The analysis presented in Table W2.5.4 reveals that differences between the semi-automated and human content are minor along most dimensions.

Table W2.5.4: Consumer Content Perception (Computational Analysis)

Dimension	Descriptives (Mean, SD) ¹				Kruskal Wallis ²			
	Revised Machine	Real SEO Experts	Quasi Experts	Novices	χ^2	η^2	df	p
Concreteness	323.10 (7.45)	326.00 (5.37)	321.30 (7.48)	318.60 (4.28)	9.67	.15	3	.021*
Familiarity	574.14 (7.95)	578.14 (12.73)	579.22 (9.33)	581.47 (9.14)	7.14	.11	3	.067
Emotionality	3.28 (.66)	3.33 (.38)	3.47 (.55)	3.53 (.47)	3.07	.05	3	.380
Emotional Valence	6.15 (.89)	6.23 (.86)	6.45 (.77)	6.69 (.72)	3.70	.06	3	.296
Negations	.004 (.005)	.005 (.003)	.006 (.003)	.007 (.006)	3.28	.05	3	.351
Interrogatives	.011 (.006)	.009 (.004)	.013 (.006)	.013 (.008)	2.31	.04	3	.509
Causation	.028 (.009)	.030 (.013)	.032 (.015)	.026 (.009)	2.07	.03	3	.558
Certainty	.011 (.005)	.013 (.005)	.021 (.009)	.019 (.009)	16.24	.25	3	.001**
Tentativeness	.022 (.010)	.027 (.014)	.022 (.010)	.022 (.009)	1.44	.02	3	.697
Differentiation	.020 (.009)	.026 (.014)	.021 (.009)	.021 (.011)	1.25	.02	3	.740
Focus on future	.009 (.006)	.013 (.006)	.011 (.006)	.015 (.007)	8.54	.13	3	.036*

¹Dimension scales: for concreteness, familiarity scale range: 100 (abstract, unfamiliar) to 700 (concrete, familiar), emotionality scale range: 0 (no emotion) to 9 (high emotion), emotional valence scale range: 0 (highly negative) to 9 (highly positive); other dimensions like negations, interrogatives, etc., represent percentages of total words in the text;

²Statistical significance codes: *0.05 level, **0.01 level; n=66;

2.6 Website Engagement

Having compared performance in terms of consumer perceptions and linguistic content, we next examine the impact of using semi-automated content on firm performance in terms of consumers' engagement with the website (e.g., Bronnenberg et al 2016, Jerath et al 2014, Edelman and Zhenyu 2016). We collect website traffic data for 412 days after the experimental content was posted. During this time, the content received 254 page views from 122 unique website visits arising from organic search results. In addition to the findings discussed in the manuscript that indicate improved performance of semi-automated content relative to human-generated content on a number of dimensions (details are reported in Table W2.6.1), we also find that the semi-automated content results in longer visits per visited page ($\chi^2(3)=167.15$, $p<.000$), suggesting better content performance (Danaher et al. 2006).

Table W2.6.1: User Behavior (Organic Search Source Only)

Dimension	Descriptives (Σ)				One-Sample Chi-Squared ¹		
	Revised Machine	Real SEO Experts	Quasi Experts	Novices	χ^2	<i>df</i>	<i>p</i>
No. of Pages with Pageviews	16	3	5	10	11.88	3	.007**
No. of Pages with Pageviews in %	84.21	33.33	26.32	52.63	40.98	3	.000**
Pageviews	172	16	18	48	257.31	3	.000**
Unique Pageviews	84	6	9	23	130.52	3	.000**
Entrances	76	6	9	21	114.21	3	.000**
Exit Rate (means)	.41	.28	.32	.36	-	-	-
Bounce Rate	.00	.00	.00	.00	-	-	-
Avg. Usage Duration (Abs., sums)	3671	262	455	473	6639.40	3	.000**
Avg. Usage Duration (Rel.) ²	229	87	91	47	167.15	3	.000**
Returning Visitors (Abs.)	88	10	9	25	127.09	3	.000**
Returning Visitors (Rel.) ²	5.50	3.33	1.80	2.50	-	-	-
Buying Affinity (Abs.) ³	4097	276	429	983	6670.60	3	.000**
Buying Affinity (Rel.) ^{2,4}	256	92	86	98	152.43	3	.000**
Exp. Sales (for U.P.*100) ⁵	168	12	18	46	151.30	3	.000**

¹Statistical significance codes: *0.05 level, **0.01 level;

²(Rel.) = the absolute value (Abs.) divided by No_of_Pages_with_Pageviews

³Buying Affinity (Abs.) = Unique_Pageviews*Willingness_to_Buy (survey measured);

⁴Buying Affinity (Rel.) = Buying_Affinity (Abs.)/No_of_Pages_with_Pageviews;

⁵Exp. Sales (for U.P.*100) = (Unique_Pageviews/100*Expected_Sales_Rate)*100, where the expected sales rate is 2% (obtained from past company reports);

Table W2.6.2 reports statistics for the user behavior for visitors coming from direct links (e.g., links in emails, on other webpages, etc.) to the focal experimental pages on the website.

Table W2.6.2: User Behavior (Direct Links Source Only)

Dimension	Descriptives (Σ)				One-Sample Chi-Squared ¹		
	Revised Machine	Real SEO Experts	Quasi Experts	Novices	χ^2	<i>df</i>	<i>p</i>
No. of Pages with Pageviews	19	9	19	19	-	-	-
No. of Pages with Pageviews in %	100	100	100	100	-	-	-
Pageviews	545	126	257	515	342.65	3	.000**
Unique Pageviews	270	65	131	257	164.05	3	.000**
Entrances	226	47	95	222	166.57	3	.000**
Exit Rate (means)	.35	.35	.35	.35	-	-	-
Bounce Rate	.04	.07	.04	.04	-	-	-
Avg. Usage Duration (Abs., sums)	705	189	536	317	361.40	3	.000**
Avg. Usage Duration (Rel.) ²	37	21	28	17	8.96	3	.029*
Returning Visitors (Abs.)	275	61	126	258	178.81	3	.000**
Returning Visitors (Rel.) ²	14.47	6.77	6.63	13.57	-	-	-
Buying Affinity (Abs.) ³	10021	3041	5846	12760	7066.10	3	.000**
Buying Affinity (Rel.) ^{2,4}	418	338	308	638	156.99	3	.000**
Exp. Sales (for U.P.*100) ⁵	540	130	262	514	328.11	3	.000**

¹Statistical significance codes: *0.05 level, **0.01 level;

²(Rel.) = the absolute value (Abs.) divided by No_of_Pages_with_Pageviews

³Buying Affinity (Abs.) = Unique_Pageviews*Willingness_to_Buy (survey measured);

⁴Buying Affinity (Rel.) = Buying_Affinity (Abs.)/No._of_Pages_with_Pageviews;

⁵Exp. Sales (for U.P.*100) = (Unique_Pageviews/100*Expected_Sales_Rate)*100, where the expected sales rate is 2% (obtained from past company reports);

2.7 Content Production Cost Calculation Details

Table 4 in the article presents the cost of content production and savings induces by applying the content writing machine. In this section, we elaborate on the calculations for Table 4. We took available working times and salary statistics for the human reviser / SEO expert necessary (i.e., as stated in the Table's footer: 39 hours available working time per week, 1,567 hours per year; 45,000 € of salary per year) and calculated the times, cost and possible output per year when using the manual way vs. the machine for text generation. To estimate the time spent

per unit of content, information was provided by both the company and the experiment participants. Based on this information, we can calculate expected outputs and labor cost per year. The calculation in Table 4 in the main manuscript are based on the following inputs, with the values from Table 4 appearing in quotations:

- “Human labor time for content production” = empirically determined
- “Server cost per unit (€)” = empirically determined
- For human groups: “Produced content units” = 1,567 hours per year / “Median (hours)”
For example, for the column “Company (real)” in Table 4: $1,567/9.5 \approx 164.95$
- For machine to keep total costs at 45,000 € (i.e., the same as the human costs): “Produced content units” = $45,000 \text{ €} / (\text{“Labour cost per unit (€)”} + \text{“Server cost per unit (€)”})$. For example, for the column “Revised Machine” in Table 4: $45,000/(15.79+5.00) \approx 2,164.03$
- “Production level (%)” = $(\text{“Produced content units” (e.g., of the revised machine)} / \text{“Produced content units” Company (real)}) * 100 - 100$. For example, for the column “Revised Machine” in Table 4: $(2,164.03/164.95) * 100 - 100 \approx 1,211.95$
- “Labor cost per unit (€)” = $45,000 \text{ €} / \text{“Produced content units”}$. For example, for the column “Company (real)” in Table 4: $45,000/164.95 \approx 272.81$
- “Cost for 164.95 units (€)” = $(\text{“Labour cost per unit (€)”} + \text{“Server cost per unit (€)”}) * 164.95$. For example, for the column “Company (real)” in Table 4:
 $272.81 + 0 * 164.95 \approx 45,000$
- “Cost for 2,164.03 units (€)” = $(\text{“Labour cost per unit (€)”} + \text{“Server cost per unit (€)”}) * 2,164.03$. For example, for the column “Company (real)” in Table 4:
 $272.81 + 0 * 2,164.03 \approx 590,369$

- “Produced content units” (“Possible real financial impact (2015 to 2019)”) = empirically determined
- “Cost (€)” = $439 * (\text{“Labour cost per unit (€)”} + \text{“Server cost per unit (€)”})$. For example, for the column “Company (real)” in Table 4: $439 * (272.81 + 0) \approx 119,765$
- “Possible savings (€)” = “Cost (€)” of the Company (Real) – “Cost (€)” of specific comparison group. For example, for the column “Revised Machine” in Table 4: $119,765 - 9,127 \approx 110,638$

3 Education Sector Field Study

In this section we present details of the event history field study we conducted in the educational sector setting as reported in the article.

3.1 Keywords to Optimize for & Keyword Statistics

Table W3.1 depicts the keywords (search queries) the experimental groups in the education sector field study produced optimal content for, including keyword statistics and descriptive statistics for the ranking performance of the revised machine content in the search engine. The basic keyword statistics reported in Table W3.1 include the average monthly search volume (i.e., rounded numbers of individuals that search on average per month for the keyword as provided by the search engine), the paid keyword competition, keyword length, etc. The keywords were selected such that they reflect target content and search queries in co-ordination with the educational institution running the experiments. The target keywords include both short and long

tail keywords, most of them are low search volume and characterized by high keyword competition.

Table W3.1: Keywords for Field Experiments

Keyword	Descriptives								
	Avg. monthly search volume	Competition	Competition index	Keyword length	Mean revised machine ranking	Median revised machine ranking	SD revised machine ranking	IQR revised machine ranking	% of days revised machine was in ranking
quantitative marketing	10	low	-	2	6.30	5	4.80	1	100.00
quantitative marketing research	10	low	0	3	8.13	6.5	6.60	5.25	93.00
quantitative marketing program	-	-	-	3	2.03	2	0.16	0	88.40
marketing research seminar series	-	-	-	4	-	-	-	-	00.00
deep learning marketing	10	low	43	3	33.70	12	58.30	2	18.60
machine learning in marketing	10	low	16	4	84.10	98	39.10	64.5	100.00
machine learning in marketing education	-	-	-	5	1.00	1	0.00	0	34.90
digital marketing and machine learning	0	-	-	5	10.70	10	1.10	1	100.00
natural language processing in marketing	10	-	-	5	10.60	11	0.70	1	95.30
artificial intelligence in marketing	50	low	32	4	36.70	35	11.60	11.5	16.30
ai in marketing	10	low	30	3	29.00	29	0.00	0	1.00
ai in digital marketing	10	mid	36	4	77.10	78	5.70	8	67.40
ai in social media marketing	10	low	14	5	74.90	83	22.20	7.75	62.80
marketing with ai	10	-	-	3	32.50	27	22.90	14.5	46.50
marketing automation	320	mid	55	2	86.00	87	7.50	7.5	14.00
customer analytics	20	low	26	2	23.00	23	4.20	3	4.65
customer segmentation with machine learning	0	-	-	5	9.50	10	1.80	2	100.00
quantitative market research methods	10	-	-	4	8.70	10	4.20	8	86.00
business analytics education	10	-	-	3	1.00	1	0.00	0	1.00
career in marketing research	10	low	0	4	7.10	7	0.40	0	32.60
career opportunities in marketing	10	low	0	4	34.00	25	15.60	13.5	6.98
methods of marketing analytics	-	-	-	4	39.00	39	0.00	0	4.65
understanding digital marketing analytics	-	-	-	4	2.00	2	0.00	0	79.10
marketing phd career opportunities	-	-	-	4	2.00	2	0.60	0	58.10
quantitative marketing phd	10	-	-	3	4.60	1	9.00	1	93.00
doctorate PHD program in marketing	-	-	-	5	15.00	13	7.60	2.5	100.00
master program in marketing	10	mid	57	4	2.30	2	0.90	0	10.00
service marketing research	10	low	0	3	4.20	1	4.10	6	100.00
research in service marketing	-	-	-	4	4.70	5	2.30	3	100.00
academic research in service marketing	-	-	-	5	1.50	1	3.20	0	100.00
marketing institute college	-	-	-	3	10.60	3	20.30	2	100.00

Entries that display “-“ mean that the search engine keyword tool did not provide specific information.

3.2 Supplemental Content Performance Tests

Table W3.2.1 reports group comparison tests for the search engine ranking performance of the experimental groups “revised machine” (printed in bold) and “human” using Wilcoxon rank sum tests. We find that the revised machine outperforms the human content generating group in terms of the number of pages that made it into the top10 search engine ranking and mean ranking performance.

Table W3.2.1: Search Engine Rankings Performance Comparison (Education Sector)

Dimension	Group	Descriptives					Wilcoxon rank sum ²			
		n _p ¹	Median	(IQR)	Min	Max	W	z	r	p
Pages in ranking / day	Revised Machine	30	18.00	(4.00)	12	22	908.00	-2.04	-.23	.041*
	Human		19.00	(2.00)	16	22				
Pages in top 10 / day	Revised Machine	30	11.50	(3.00)	7.00	14.00	24.50	7.20	.82	.000**
	Human		5.00	(2.00)	3.00	9.00				
Mean rankings / day	Revised Machine	30	17.57	(9.44)	5.60	30.12	1270.5	-5.79	-.66	.000**
	Human		26.22	(2.72)	19.44	30.30				

¹n_p=number of pages per experimental group. n=77 (days); ²Two-tailed tests; statistical significance codes: *0.05 level, **0.01 level; Compared numbers are daily aggregate numbers.

Table W3.2.2 reports quality score statistics and significance testing results for the experimental groups (i.e., revised machine versus humans), the top 10 and the lowest ranked pages on the five

quality score components. The findings are consistent with our findings from the IT service industry experiment.

Table W3.2.2: Quality Score Components Group Comparisons with
Top 10 Ranked Websites (Education Sector)

Quality Score Component	Group	Descriptives				Wilcoxon rank sum ¹			
		Median (IQR)		Min	Max	W	z	r	p
Topic (<i>s_a</i>)	Top 10	.44	(.20)	.15	.63				
	Revised machine	.57	(.13)	.29	.76	205	3.72	.48	.000**
	Raw machine	.56	(.12)	.25	.76	223	3.42	.44	.000**
	Humans	.43	(.09)	.25	.64	476	-.37	-.05	.708
	Worst 10	.20	(.08)	.12	.38	803	-6.18	-.81	.000**
Keywords (<i>s_k</i>)	Top 10	.47	(.23)	.09	.70				
	Revised machine	.65	(.16)	.27	.85	194	3.90	.50	.000**
	Raw machine	.63	(.15)	.17	.84	190	3.97	.51	.000**
	Humans	.49	(.16)	.25	.73	393	.83	.11	.406
	Worst 10	.15	(.11)	.06	.42	805	-6.23	-.81	.000**
Uniqueness (<i>s_a</i>)	Top 10	.95	(.06)	.78	.99				
	Revised Machine	.94	(.08)	.73	1.00	527	-1.13	-.15	.258
	Raw machine	.90	(.09)	.72	.99	623	-2.57	-.33	.010*
	Humans	.99	(.05)	.07	1.00	234	3.19	.41	.001**
	Worst 10	.94	(.03)	.86	.98	459	.35	.05	.724
Readability similarity (<i>s_r</i>)	Top 10	.57	(.07)	.46	1.00				
	Revised Machine	.73	(.25)	.47	1.00	135.5	4.64	.60	.000**
	Raw machine	.78	(.23)	.45	1.00	85	5.39	.70	.000**
	Humans	.47	(.54)	.00	1.00	544.5	-1.39	-.18	.165
	Worst 10	.41	(.23)	.15	.70	726	-4.41	-.57	.000**
Naturalness similarity (<i>s_n</i>)	Top 10	.57	(.08)	.44	.75				
	Revised Machine	.54	(.50)	.17	1.00	461	-0.15	-.02	.876
	Raw machine	.67	(.50)	.33	1.00	393	.84	.11	.403
	Humans	.17	(.25)	.00	1.00	780	-4.89	-.63	.000**
	Worst 10	.41	(.16)	.20	.88	725.5	-4.40	-.57	.000**

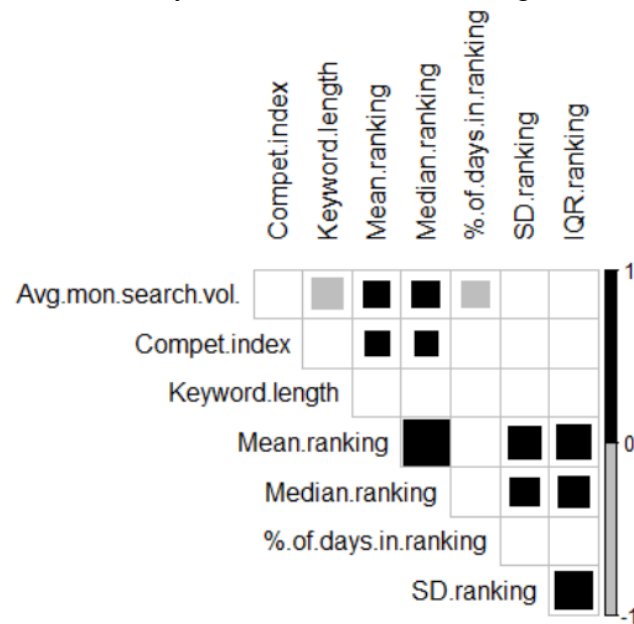
¹ Two-tailed tests; statistical significance codes: *0.05 level, **0.01 level;

4 Keyword Boundary Conditions

In this section, we investigate possible keyword boundary conditions (e.g., keyword length, competition levels, etc.) on the revised content machine search engine ranking performance using the data and keywords reported in Tables W2.1 (IT service sector) and W3.1 (education sector).

Figure W4 contains a series of Kendall's tau b correlations in which we try to ascertain if factors such as keyword length, search volume, or competition have an effect on the search engine ranking performance of the revised machine content. Figure W4 shows that the machine-generated content does not perform as well (i.e., higher values of mean and median search engine ranking positions) for more popular search terms (higher average monthly search volume and more competition). In addition, a higher average monthly search volume is also associated with a lower percentage of days for which machine-generated content is ranked in the search engine over our observation period. The keyword length does not play an important role.

Figure W4: Keyword Boundary Conditions on Search Engine Ranking Performance



Kendall's tau b correlations: ■ = positive correlation, ■ = negative correlation, □ = statistically non significant correlation (at a 0.05 alpha level), a bigger square represents a higher correlation coefficient;

5 Appendix References

- Baayen RH, Shafaei-Bajestan E (2019) Analyzing linguistic data: A practical introduction to statistics. Package ‘languageR’. Version 1.5.0. CRAN. Accessed May 20, 2019, <https://cran.r-project.org/web/packages/languageR/languageR.pdf>
- Benoit K, Watanabe K, Wang H, Nulty P, Obeng A, Müller S, Matsuo A, (2018) “quanteda: An R package for the quantitative analysis of textual data.” *Journal of Open Source Software*. 3(30). <https://doi.org/10.21105/joss.00774>
- Berger J, Sherman G, Ungar L (2020b) TextAnalyzer. Accessed November 11, 2020, <http://textanalyzer.org>
- Bronnenberg BJ, Kim JB, Mela CF (2016) Zooming in on choice: How do consumers search for cameras online? *Marketing Science*. 35(5):693-712.
- Danaher PJ, Mullarkey GW, Essegai S (2006) Factors affecting website visit duration: A cross-domain analysis. *Journal of Marketing Research*. 43(2):182-194.
- Edelman B, Zhenyu L (2016) Design of search engine services: Channel interdependence in search engine results. *Journal of Marketing Research*. 53(6):881-900.
- Flanigan, AJ, Metzger, MJ (2007) The role of site features, user attributes, and information verification behaviors on the perceived credibility of web-based information. *New Media & Society*. 9(2):319-342. <https://doi.org/10.1177/1461444807075015>
- Jerath K, Ma L, Park YH (2014) Consumer click behavior at a search engine: The role of keyword popularity. *Journal of Marketing Research*. 51(4):480-486.
- Kamoen N, Holleman B, Bergh H (2013) Positive, negative, and bipolar questions: The effect of question polarity on ratings of text readability. *Survey Research Methods*. 7(3):181-189.
- Liu J, Toubia O (2018) A semantic approach for estimating consumer content preferences from online search queries. *Marketing Science*. 37(6):930-952.
- Maechler M, Rousseeuw P, Croux C, Todorov V, Ruckstuhl A, Salibian-Barrera M, Verbeke T, Koller M, Conceicao ELT, Palma MA (2020) Basic robust statistics. Package ‘robustbase’. Version 0.93-6. CRAN. Accessed May 20, 2020, <https://cran.r-project.org/web/packages/robustbase/robustbase.pdf>
- Pennebaker JW, Booth RJ, Boyd RL, Francis ME (2015) Linguistic inquiry and word count: LIWC2015. Austin, TX: Pennebaker Conglomerates. Accessed November 1, 2020, www.LIWC.net.

Pitler E, Nenkova A (2008) Revisiting Readability: A unified framework for predicting text quality. *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*. 186-195.

Radford A, Narasimhan K, Salimans T, Sutskever I (2018) Improving language understanding by generative pre-training. OpenAI.

Roberts C (2010) Correlations among variables in message and messenger credibility scales. *American Behavioral Scientist*. 54(1):43-56.

Rocklage MD, Rucker DD, Nordgren LF (2018) Persuasion, emotion and language: the intent to persuade transforms language via emotionality. *Psychological Science*. 29(5):749-760.

Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser L, Polosukhin I (2017) Attention is all you need. *31st Conference on Neural Information Processing Systems (NIPS 2017)*. 1-15.