

The Art of Inhibiting Patterns: Gaussian Process Optimization of Hyphenation

Anonymous ACL submission

Abstract

The quality of hyphenation matters. It does matter when rendering texts on small mobile screens, when rendering HTML in any browser, and it matters in the narrow columns of academic papers, such as in this abstract. Having an automated, robust, efficient, and generalizable solution for all languages in use would have saved billions of hours of scrolling and millions of tons of paper.

The half-century of attempts to solve the hyphenation problem yielded usable solutions with $\text{T}_\text{E}\text{X}$ hyphenation patterns, but with low coverage and errors caused mainly by insufficient exception handling, compound words, and language changes (including new words, transcriptions). The setting of parameters of the pattern generation process is still done by trial and error.

We have designed a new method of Gaussian process optimization of hyphenation pattern generation by the `patgen` program. It optimizes threshold setting for coverage and inhibition levels of pattern generation. We utilized a collection of 18 diverse datasets comprising a hyphenated word list for 15 languages and generated nearly optimal hyphenation patterns with respect to precision, recall, and generalization (10-fold cross-validation) optimization criteria. The $F_{1/7}$ score reaches 0.9992 for consistently marked word lists, such as Czech and Slovak.

A newly designed optimization technique for the pattern generation methodology has set new quality standards, approaching an optimal solution. Its deployment in typesetting systems like $\text{T}_\text{E}\text{X}$, InDesign, or in text processors, as well as in CSS hyphenation in web browsers, would yield a huge impact.

“Perfection is not attainable, but if we chase perfection, we can catch excellence.” Vince Lombardi

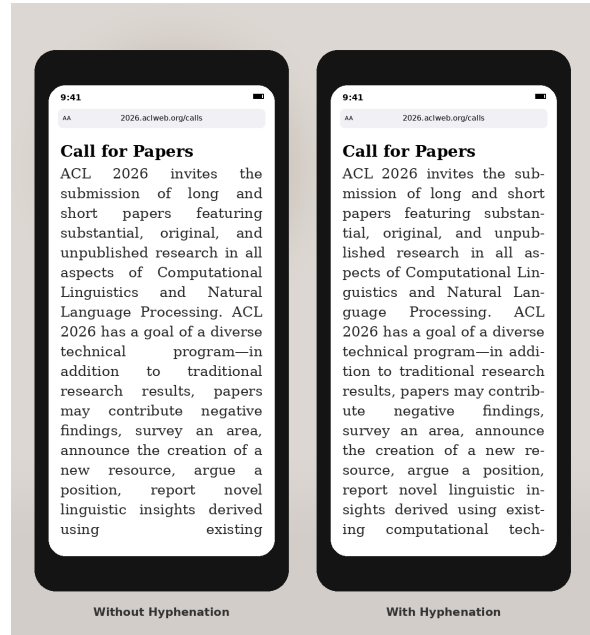


Figure 1: Comparison of text rendering on a mobile device. Note the 'rivers' of whitespace without hyphenation (left) and the improved spacing density with hyphenation and fewer lines for scrolling (right).

1 Motivation

The quality of hyphenation matters. It matters when rendering HTML texts on small mobile screens, as shown in Figure 1. It matters in the narrow columns of academic papers, such as in the two-column papers in ACL Anthology, as in this paragraph.

Having an automated, robust, efficient, and generalizable solution for all languages in use would have saved billions of hours of scrolling and millions of tons of paper.

2 The Development of Hyphenation Algorithms

Approaches have evolved from pattern-based and rule-based systems to various statistical machine learning methods, including neural.

058 **Pattern-Based Hyphenation**

059 The digital pursuit of hyphenation quality began in
060 1983 when Donald Knuth’s doctoral student Frank
061 Liang (Liang, 1983) designed a supervised learn-
062 ing method to generate hyphenation *patterns* from
063 a list of already hyphenated words of Webster’s
064 Pocket Dictionary, for use in T_EX (Knuth, 1986,
065 Appendix H). The algorithm implemented in the
066 patgen program (Liang and Breitenlohner, 2020)
067 employs a dynamic packed trie data structure to
068 efficiently store patterns and identify hyphenation
069 points based on local *context* windows. It oper-
070 ates on a multilevel system where patterns can ei-
071 ther enable or inhibit the hyphenation point, based
072 on weighted values. This method is renowned for
073 its exceptional space efficiency and high coverage.
074 Liang’s hyphen.tex with 4,447 patterns covers
075 89.3% of hyphens in the input word list.

076 Sojka (1995) showed that the patterns could be
077 generated with coverage of 98.37% for German,
078 and that the method could be used for segmenta-
079 tion of compounds, too.

080 **Finite State Methods** Bouma (2003) explores
081 the formalization of hyphenation using Finite State
082 Transducers (FSTs). Unlike the specialized packed
083 trie implementation found in T_EX, this approach
084 utilizes standard finite state calculus, where hy-
085 phenation patterns and rules are expressed as regu-
086 lar expressions. This method demonstrates that hy-
087 phenation can be effectively modeled using deter-
088 ministic automata, providing a rigorous mathemat-
089 ical framework that allows the seamless composi-
090 tion of hyphenation with other finite-state morpho-
091 logical processes while maintaining accuracy com-
092 parable (99.1% accuracy) to traditional pattern-
093 based systems (99.8% accuracy) on running Dutch
094 text (Bouma, 2003, Table 5).

095 **Subword Segmentation by Transformers** While
096 distinct from typographic hyphenation (which fo-
097 cuses on line breaking), subword segmentation is
098 the dominant paradigm in modern NLP for train-
099 ing Large Language Models (LLM). Sennrich et
100 al. (Sennrich et al., 2016) introduced Byte-Pair
101 Encoding (BPE) for neural machine translation,
102 which iteratively merges frequent character pairs
103 to form subword units.

104 **Language-Agnostic Neural Syllabification** A
105 language-agnostic approach, framing syllabifica-
106 tion as a sequence labeling task utilizing Bidirec-
107 tional Long Short-Term Memory (Bi-LSTM) net-

works, has been proposed by Krantz et al. (2019).
By treating hyphenation points as classification
targets between characters, they demonstrate that
high-accuracy segmentation can be achieved di-
rectly from raw text data without reliance on
language-specific rules or manual feature engineer-
ing. This contrasts with the patgen approach
by prioritizing generalization capabilities through
deep learning over the storage efficiency of packed
tries. Franek (Franěk, 2025) measures trade-offs
between finite-state and minimal LLM-based ap-
proaches.

The most widespread method for hyphenation
is pattern-based hyphenation, due to its high accu-
racy and generalization (Sojka and Sojka, 2023),
efficiency (Sojka and Sojka, 2019a), small size,
and potential applicability to any language with a
hyphenated word list.

LLMs are considered overkill (computational
waste) for simple rendering tasks on mobile de-
vices/browsers, where a footprint of less than
50 KB is sufficient with pattern-based approaches.

Even though the best results for ever-developing
languages and their new words are obtained by pat-
terns generated by patgen from the large-scale
word list of hyphenated words, a substantial por-
tion of the 90 languages in the repository *hyph-*
utf8 rely on manually created rules. To develop
hyphenation patterns, one needs to create a big,
consistently hyphenated word list, together with
heuristically optimized parameters for the patgen
generation.

This paper describes a new flexible approach
to the optimization of pattern generation that is
based on Gaussian Process (*GP*) (Rasmussen and
Williams, 2005). We start with the description of
available word lists in Section 3. We proceed with
the pattern generation process principles in Sec-
tion 4. We describe the evaluation process and met-
rics used in Section 5. The design of our optimiza-
tion by *GP* is explained, and the main results are
described in Section 6. We conclude with the sum-
mary and future work in Section 7.

3 Datasets: Hyphenated Word Lists

Metelka and Sojka (2025) recently reported hy-
phenation benchmark with hyphenated word lists.

Table 1 shows the diversity of their collection:
there are 18 word lists for 15 languages in several
alphabets (Latin, Cyrillic, and Thai). There is sig-
nificant variability between datasets (word lists) in

158	terms of their size, as well as average line lengths	208
159	and density of hyphenation marks.	209
160	They stress that the quality of the data varies	210
161	across datasets. Some of their datasets are curated	
162	and very consistent, such as <i>de/wortliste</i> (Lemberg,	
163	2025) or <i>cssh/cshyphen</i> (Sojka and Sojka, 2021).	
164	On the other hand, there are no guarantees of qual-	212
165	ity in <i>wiktionary</i> datasets, as anyone can contribute	213
166	to a Wiktionary entry, and these entries are typi-	214
167	cally not subject to professional review The dataset	215
168	diversity will contribute to the evaluation of the	216
169	new optimized pattern generation in the next sec-	217
170	tion.	218
171	4 Pattern Generation with patgen	
172	The patgen pattern generation tool (Liang, 1983)	220
173	is part of the T _E X Live 2025 distribution. It creates	221
174	patterns from hyphenated word lists that are subse-	222
175	quently used in the T _E X82 hyphenation algorithm.	223
176	For efficient lookup, the algorithm uses a dynamic	224
177	packed trie data structure for patterns.	225
178	A single run of the program consists of several	226
179	phases. First, the <i>translate file</i> is read if present,	
180	and the characters (byte sequences) it contains are	
181	mapped to patgen’s internal codes. Then the input	
182	patterns are loaded. Subsequently, the algorithm it-	
183	erates over the word list and decides which of the	
184	candidate patterns to include in the output collec-	
185	tion. Finally, the selection is evaluated on the in-	
186	put word list to obtain values <i>good</i> (hyphenation	
187	point correctly identified), <i>bad</i> (hyphenation point	
188	identified on a non-hyphenated position of a word),	
189	and <i>missed</i> (hyphenation point not identified, but a	
190	word is hyphenated on that position). Speaking in	
191	statistical terms, these correspond to true positives,	
192	false positives, and false negatives, respectively.	
193	The quality of patterns generated by patgen is	
194	influenced by the choice of algorithm parameters.	
195	We provide an overview in Table 2. It is possible to	
196	control the minimal distance of hyphenation points	
197	from the edge of the word (<code>\lefthyphenmin</code> ,	
198	<code>\righthyphenmin</code>) and lengths of the patterns ex-	
199	plored in the current run (<code>pat_start</code> , <code>pat_finish</code>).	
200	The power and the meaning of the generated pat-	
201	terns can be modified by specifying the hyphen-	
202	ation level (<code>hyph_start</code> , <code>hyph_finish</code>): patterns on	
203	higher levels overwrite those on lower levels dur-	
204	ing hyphenation; odd levels create hyphenating	
205	patterns, while even levels create inhibiting pat-	
206	terns (no hyphenation allowed). The parameters	
207	<i>good_wt</i> , <i>bad_wt</i> , and <i>thresh</i> play a key role in pat-	
	tern selection. For a pattern to be accepted, the	208
	following inequality must hold (Liang and Breiten-	209
	lohner, 2020):	210
	$\alpha * good_wt - \beta * bad_wt \geq thresh, \quad (1)$	211
	where α denotes the number of positive exam-	212
	ples (word split correctly by hyphenating pattern	213
	/ previously wrong hyphenation fixed by inhibiting	214
	pattern), and β the number of negative examples	215
	(word split incorrectly by hyphenating pattern / pre-	216
	viously correct hyphenation broken by inhibiting	217
	pattern) in the word list.	218
	5 Pattern Generation Evaluation Metrics	219
	We define the metrics for the evaluation of the hy-	220
	phenation algorithm, pattern generation process,	221
	and <i>parameter profiles</i> that serve for advanced	222
	search techniques of parameter generation. When-	223
	ever possible, we opted for the same or similar met-	224
	rics as used in prior work (Sojka and Sojka, 2023;	225
	Metelka and Sojka, 2025):	226
	Precision P: the ratio of correct hyphenations out	227
	of all identified hyphenations, defined as $P =$	228
	$\frac{good}{good+bad}$. The values <i>good</i> , resp. <i>bad</i> , and	229
	<i>missed</i> are the amounts of correctly resp. in-	230
	correctly identified hyphenation points. This	231
	is the <i>primary metric</i> , as we do not want errors	232
	to appear in the automatically typeset docu-	233
	ments.	234
	Recall R: the ratio of correctly identified hyphen-	235
	ations out of all hyphenations in the data, de-	236
	finied as $R = \frac{good}{good+missed}$. This metric, also	237
	called Coverage , is the <i>secondary metric</i> , as	238
	allowing more hyphenation points may result	239
	in more uniform interword spacing and eco-	240
	nomical savings due to shorter documents.	241
	F-score F: weighted harmonic mean of P and R	242
	(Rijsbergen, 1979). We prioritize P much	243
	more than R and use $F_{1/7}$ -score, defined as	244
	$F_{1/7} = \frac{(1+1/7)^2 * P * R}{(1/7)^2 * P + R}$. The $F_{1/7}$ -score thus	245
	serves as a joint measure of Precision and Re-	246
	call, emphasizing the greater importance of	247
	not making errors over coverage.	248
	Size: The number of trie nodes used to store the	249
	generated pattern is our <i>tertiary metric</i> . Trie	250
	depth limits the $O(1)$ time to find hyphen-	251
	ation points for a given word from patterns.	252

Table 1: Dataset collection overview: full size in kBytes for UTF-8 encoded files (*Size*), average word (segment) length in characters (*#Lines*), average number of hyphenation points per word (*#Hyphs*), `\lefthyphenmin` and `\righthyphenmin` parameters (*LRhyp*), and the number of distinct lowercase characters (*#Chars*) occurring.

<i>Language/Name</i>	<i>Size</i>	<i>#Lines</i>	<i>#Lines</i>	<i>#Hyphs</i>	<i>#Chars</i>
<i>cssk/cshyphen</i>	10,646.4	835,145	11.10	2.54	50
<i>cs/cshyphen_cstenten</i>	8,202.6	597,170	11.83	2.79	48
<i>cs/cshyphen_ujc</i>	1,377.1	104,277	11.48	2.59	48
<i>cs/wiktionary</i>	632.0	53,962	9.95	2.16	48
<i>de/wiktionary</i>	12,821.6	807,813	15.01	2.96	86
<i>de/wortliste</i>	10,035.4	604,970	15.78	2.99	47
<i>el/wiktionary</i>	603.9	28,208	11.90	2.87	52
<i>es/wiktionary</i>	11,753.4	816,723	13.36	3.27	50
<i>is/hyphenation-is</i>	3,037.8	218,308	12.12	1.91	45
<i>it/wiktionary</i>	981.1	77,351	11.96	2.90	39
<i>ms/wiktionary</i>	13.6	1,247	8.52	1.81	69
<i>nl/wiktionary</i>	8,124.5	528,598	14.73	2.82	68
<i>pl/wiktionary</i>	1,713.6	141,916	10.75	2.20	52
<i>pt/wiktionary</i>	541.6	44,321	11.17	2.74	40
<i>ru/wiktionary</i>	3,795.6	152,608	13.10	2.74	51
<i>th/orchid</i>	3,567.1	34,868	40.69	6.02	95
<i>tr/wiktionary</i>	187.2	15,619	10.41	2.37	42
<i>uk/wiktionary</i>	351.6	15,701	11.79	2.65	41

Table 2: Overview of patgen parameters in pattern generation profiles

Name	Value range	Description
<code>\lefthyphenmin</code>	1 – 99	Length of word prefix without hyphenation, typically 1
<code>\righthyphenmin</code>	1 – 99	Length of word suffix without hyphenation, typically 1 or 2
<i>hyph_start</i>	1 – 9	Initial hyphenation level, typically 1
<i>hyph_finish</i>	1 – 9	Final hyphenation level, typically 4 or 5
<i>pat_start, pat_finish</i>	1 – 15	Minimal resp. maximal pattern length in given level, ranges typically increase in higher levels
<i>good_wt</i>	integer	Weight of examples in favor of the pattern
<i>bad_wt</i>	integer	Weight of examples against the pattern
<i>thresh</i>	integer	Evaluation threshold

Generalization: To assess the performance of the created patterns on unseen data, a 10-fold cross-validation is performed: We generate patterns from a train split of the data, measure their $F_{1/7}$ on the unseen test split, and average the values across the folds.

Parameter Profiles To compare different strategies of pattern generation, we define *parameter profile* as a set of adjustable parameter values that fully determine the patgen run, e.g., the quality of the resulting patterns. This means specifying pattern length window (*pat_start, pat_finish*) and

weights (*good_wt, bad_wt, thresh*) for every hyphenation level as shown in Table 2.

Over the past three decades of experimentation with pattern generation, numerous profiles for various word lists have been developed. In Table 3, we show three profiles that gave the best results in prior work. One is adapted from earlier works (Sojka and Sojka, 2021, Table 2, *correct optimized* profile) *cshyphen*, and two from (Lemberg, 2025) *wortliste8*, and *wortliste* that contain the first 4 levels of the earlier.

Each profile defines a point p in patgen n -dimensional hyperparameter space, which sets the

Table 3: Parameter profile examples found by decades of empirical investigations. Parameters *pat_start* and *pat_finish* are combined into column “Lengths” and *good_wt*, *bad_wt*, and *thresh* into column “Params”.

Profile	Level	Lengths	Params
<i>cshyphen</i>	1	1 3	1 5 1
	2	1 3	1 5 1
	3	2 6	1 3 1
	4	2 7	1 3 1
<i>wortliste</i>	1	1 3	2 3 1
	2	2 4	1 5 1
	3	3 5	1 6 1
	4	4 6	1 7 1

Profile	Level	Lengths	Params
<i>wortliste8</i>	1	1 3	2 3 1
	2	2 4	1 5 1
	3	3 5	1 6 1
	4	4 6	1 7 1
	5	5 12	1 8 1
	6	6 12	1 9 1
	7	7 12	1 9 1
	8	8 12	1 9 1

patterns’ quality and performance. Hand tuning of profiles for each word list is very laborious and time-consuming work for patgen experts. A natural research question arises: can we find an *algorithmic way* better than trial and error to compute an optimal profile p , according to some objective function?

6 Optimal Pattern Generation with Gaussian Processes

The core idea of this paper is to define and use a Gaussian Process (Rasmussen and Williams, 2005; Snoek et al., 2012) to compute a profile p algorithmically in such a way that our three metrics are *jointly* optimized.

Formally, a Gaussian Process is specified by two functions: mean and covariance functions. We define our mean function m that represents the expected value of profile p as

$$m(p) = F_{1/7} - \text{trie_penalty}. \quad (2)$$

We are looking for a profile p that maximizes $m(p)$ so that it is as close to 1 as possible. The rationale behind this design is as follows.

The first term, $F_{1/7}$, maximizes both primary and secondary metrics (precision and coverage). The preference for precision over recall is proper because false positives (incorrect hyphenations) are more severe than false negatives (missed hyphenations), which is typically the case for hyphenation problems. We followed the practice of prior work with this tradeoff setting, where bad hyphenation is 49 times worse than missed hyphenation.

The second term optimizes our tertiary metrics, pattern size, as it encourages smaller, more generalizable pattern sets. The trie node count correlates with the generalization. Smaller tries with shorter

patterns and shorter contexts often generalize better to unseen words.

$$\text{trie_penalty} = \text{trie_weight} \frac{\text{trie_nodes}}{\text{trie_normalizer}} \quad (3)$$

The second critical component of a *GP* is the covariance function (kernel) that defines the relationship between input points p and p' . If the points are close to each other, their output values $m(p)$ and $m(p')$ should also be close. Our implementation of *GP* optimizer finds the best profile (and patterns) within n -dimensional parameter space. It uses GaussianProcessRegressor with Matern and WhiteKernel kernels from sklearn.Gaussian_process. It uses Upper Confidence Bound $UCB(x) = \mu(x) + \kappa \cdot \sigma(x)$ acquisition for exploration-exploitation balance.

Results We computed the best profiles using *GP* and compared them with the current State-of-the-Art (SOTA) profiles *cshyphen* and *wortliste*. We used the following optimization settings named *GPopt4(p)*:

objective: $F_{1/7}$ with trie size penalty as defined in Eq. (2).

number of generation levels: As there were no published examples of languages that would require a deeper hierarchy of exceptions than two inhibiting levels, we run *GP* for 4 levels. In each level, we have to decide parameters of Equation (1): *good_wt*, *bad_wt*, and *thresh*.

parameter ranges: To minimize the search space for 4 equations 1, we decided to force the same *thresh* for each level, bounded by $[1, 5]$. We search for *bad_wt* in $[1, 30]$, with fixed *good_wt* = 3.

346	parameter dimension n: $n = 5$, <i>bad_wt</i> for four	Summary statistics <i>GPopt4(p)</i> profile resulted	386
347	levels and one shared <i>thresh</i> .	in	387
348	pattern lengths ranges: We start with short pat-	$F_{1/7}$ improved for 11/18 datasets: <i>cssh/cshyphen</i> ,	388
349	terns to continue with a longer context for in-	<i>cs/cshyphen_cstenten</i> , <i>de/wiktionary</i> ,	389
350	hibition. Range pairs (<i>min_length</i> , <i>max_length</i>)	<i>el/wiktionary</i> , <i>it/wiktionary</i> , <i>ms/wiktionary</i> ,	390
351	for candidate patterns in each level are:	<i>nl/wiktionary</i> , <i>pl/wiktionary</i> , <i>pt/wiktionary</i> ,	391
352	$[(1,4), (2,5), (2,6), (2,7)]$.	<i>th/wiktionary</i> , <i>uk/wiktionary</i>	392
353	trie_weight: 0.0005 as used in Eq. (3).	pattern tries smaller than best previous baseline	393
354	trie_normalizer: 25000 as used in Eq. (3);	for 7/18 datasets: <i>cssh/cshyphen</i> ,	394
355	widespread natural language uses at least	<i>el/wiktionary</i> , <i>is/wiktionary</i> , <i>ms/wiktionary</i> ,	395
356	25000 lemmas.	<i>nl/wiktionary</i> , <i>uk/wiktionary</i> , <i>es/wiktionary</i>	396
357	iterations: 100 at batch size 1.	Efficiency and effectiveness remarks The en-	397
358	UCB kappa: $\kappa = 2.5$ seems to be the right bal-	tire optimization of Gaussian processes for all 18	398
359	ance between exploration and exploitation.	word lists, along with follow-up cross-validation	399
360	To assess the generalization properties of the cre-	computations, could be rerun on a commodity lap-	400
361	ated patterns, a 10-fold cross-validation was per-	top without a GPU in under a day.	401
362	formed. We generate patterns from a train split of	While neural models such as Transformer/Bi-	402
363	the data, measure their performance on the unseen	LSTM hyphenator may be pushed to achieve	403
364	test split, and average the values across the folds.	99.9% accuracy, they require 1000 – 10000× the	404
365	We utilize the diverse dataset introduced by	memory.	405
366	Metelka and Sojka (2025) . In Table 4 we report	Our method achieves a higher 99.999% accu-	406
367	the results of our experiments run on this dataset.	racy with a 20 KB footprint, making it the only vi-	407
368	We have prepared several visualizations clarify-	able solution for embedding in mobile OS kernels	408
369	ing our <i>GP</i> method in Appendix A.	or low-power e-readers. The inference speed is ap-	409
370	Key success The primary dataset <i>cssh/GPopt4(p)</i>	proximately 1000× faster when measured by the	410
371	for the paper (Czech+Slovak word list combined)	number of hyphenated words per second.	411
372	shows that <i>GP</i> beats <i>all</i> baselines:	Reproducibility All reported experiments are	412
373	• <i>GPopt4(p)</i> : $F_{1/7} = 0.9992$, and trie size =	fully reproducible. Command	413
374	18888.	<code>python -m scripts.optimize --lang cssh/cshyphen \</code>	414
375	• Best baseline: $F_{1/7} = 0.9986$, and trie size =	<code>--iterations 100 --batch-size 5 \</code>	415
376	25302 (<i>cshyphen</i> profile).	<code>--objective f17_trie --good-weight 3 \</code>	416
377	• Improvement: +0.0006 on $F_{1/7}$ (higher pre-	<code>--max-bad-weight 30 --max-threshold 1 \</code>	417
378	cision), and –6414 trie size (–25% smaller).	<code>--ucb-kappa 2.5 --trie-weight 0.0005 \</code>	418
379	Trie-$F_{1/7}$ tradeoff Most <i>GP</i> solutions achieve	<code>--trie-normalizer 25000</code>	419
380	higher training $F_{1/7}$ but larger tries that don't gen-	does the <i>GP</i> optimization part on the Czecho-	420
381	eralize as well.	Slovak word list.	421
382	Overfitting tendency The wide search space	All reported experiments (software and data)	422
383	(<i>bad_wt</i> in [1,30]) leads to overfitting for some	are fully reproducible and available in the anony-	423
384	datasets. For production use, one might consider	mous repository https://anonymous.4open.science/r/hyph-bench-5F86/README.md .	424
385	using tighter bounds (1–9) in matching profiles.	“If I waited for perfection...I would never write a word.”	426
		– Margaret Atwood	427
		7 Conclusion and Future Work	428
		We have designed a Gaussian Process optimized	429
		for hyphenated patterns generated by the <i>patgen</i>	430
		program.	431
		We have demonstrated that automated <i>GP</i> op-	432
		timization can find solutions that further improve	433
		<i>patgen</i> generation profiles hand-tuned for four	434

Table 4: Cross-validation results overview: $F_{1/7}$ -score (rounded to 4 decimal points) and the number of nodes in the pattern-storing trie (rounded to the lowest higher integer) for *cshyphen*, *wortliste*, and *GPopt4(p)* parameter profiles, averaged over the 10 cross-validated folds. The 5 parameters optimized by *GPopt4(p)* are 4 *bad_wt* for each level and *thresh*.

<i>Language/Name</i>	<i>cshyphen</i>		<i>wortliste</i>		<i>GPopt4(p)</i>						
	$F_{1/7}$	trie size	$F_{1/7}$	trie size	$F_{1/7}$	trie size	5 params				
<i>cssk/cshyphen</i>	0.9986	25302	0.9981	21685	0.9992	18888	2	2	6	12	1
<i>cs/cshyphen_cstenten</i>	0.9444	12052	0.9243	12121	0.9586	11929	3	1	6	7	1
<i>cs/cshyphen_ujc</i>	0.9830	17046	0.9803	14736	0.9820	26653	23	1	30	16	1
<i>cs/wiktionary</i>	0.9759	10643	0.9744	8330	0.9742	17574	27	2	29	25	1
<i>de/wiktionary</i>	0.9983	32394	0.9977	29449	0.9985	45653	20	2	22	1	1
<i>de/wortliste</i>	0.9775	47544	0.8985	43214	0.9452	47397	2	3	5	1	1
<i>el/wiktionary</i>	0.9875	3602	0.9866	3223	0.9880	3524	4	2	30	1	1
<i>es/wiktionary</i>	0.9998	2023	0.9998	1784	0.9998	1902	1	1	23	1	1
<i>is/hyphenation-is</i>	0.9839	32284	0.9802	28664	0.9825	27324	1	5	1	1	1
<i>it/wiktionary</i>	0.9970	3593	0.9968	2875	0.9971	4158	30	1	30	21	1
<i>ms/wiktionary</i>	0.9296	796	0.9244	721	0.9302	937	15	1	1	11	1
<i>nl/wiktionary</i>	0.9972	30859	0.9963	27514	0.9972	31798	3	2	10	1	1
<i>pl/wiktionary</i>	0.9913	8564	0.9904	7262	0.9915	15062	30	1	24	1	1
<i>pt/wiktionary</i>	0.9969	1761	0.9970	1422	0.9971	1925	13	1	23	1	1
<i>ru/wiktionary</i>	0.9270	53820	0.9212	42164	0.9252	101559	30	8	29	1	1
<i>th/orchid</i>	0.9622	28795	0.9488	24258	0.9681	38955	30	6	30	1	1
<i>tr/wiktionary</i>	0.9920	1406	0.9916	1273	0.9915	2064	30	1	24	1	1
<i>uk/wiktionary</i>	0.9602	4741	0.9564	4272	0.9615	4514	3	1	3	11	1

435 decades. With consistently marked and curated hy-
 436 phenated word lists, the resulting solution might be
 437 very, very close to the optimal one.

438 These results set the floor for improvements and
 439 savings in all typesetting systems and rendering en-
 440 gines that hyphenate language texts.

441 We hope that the new method will stimulate the
 442 research and development of new word lists or fur-
 443 ther pattern generation strategies. The reported re-
 444 sults for *cssk/GPopt4(p)* are encouraging for the
 445 development of patterns for sets of languages that
 446 share phonetic rules of syllabification and hyphen-
 447 ation, e.g., for universal Slavic hyphenation.

448 In future work, to achieve the above vision, the
 449 next steps are:

450 **further GP profile development** Pattern genera-
 451 tion could still be further improved (speed)
 452 by *annealing* techniques (Kirkpatrick et al.,
 453 1983; Ingber, 1989).

454 **word list development** Word lists could be ex-
 455 tended from available sources with the
 456 techniques of bootstrapping (Sojka and
 457 Sojka, 2021, 2019b) and transfer learning
 458 from (Sojka, 2025).

pattern generation for big alphabets Universal
 hyphenation and segmentation preparation
 (Shao et al., 2018) will need a new version
 of the *patgen* program capable of working
 with Unicode-size alphabets.

464 Limitations

465 The current optimization parameters of the Gaus-
 466 sian process need further investigation to minimize
 467 the heuristic setting of parameter ranges. One pos-
 468 sibility is to use wider parameter ranges with some
 469 of the efficient annealing techniques, e.g., very fast
 470 simulated re-annealing by Ingber (1989). We are
 471 aware that the current set of hyphenated word lists
 472 is limited, as is the size and consistency of data
 473 from Wiktionary. For universal use in any natu-
 474 ral language, the *GP* should be tested on a bigger
 475 database of hyphenated words.

476 While we utilized AI assistance (Claude 4.5,
 477 Gemini 3 Pro) for copy-editing and code checking
 478 in accordance with the ACL 2026 Policy on Publi-
 479 cation Ethics, all scientific claims and methodolog-
 480 ical designs were verified by the human authors.

Table 5: Comparison of profiles: *GP* optimized *GPop4(p)* versus previous best base SOTA profiles *cshyphen* or *wortliste*

<i>Language/Name</i>	<i>GP</i> $F_{1/7}$	<i>GP</i> trie $F_{1/7}$	Best base $F_{1/7}$	Best base profile	Base trie	$F_{1/7}$ better?
<i>cssk/cshyphen</i>	0.9992	18888	0.9986	cshyphen	25302	YES (+0.0006)
<i>cs/cshyphen_cstenten</i>	0.9586	15278	0.9444	cshyphen	12052	YES (+0.0142)
<i>cs/cshyphen_ujc</i>	0.9820	26653	0.9830	cshyphen	17046	NO (−0.0010)
<i>cs/wiktionary</i>	0.9742	17574	0.9759	cshyphen	10643	NO (−0.0017)
<i>de/wiktionary</i>	0.9985	45653	0.9983	cshyphen	32394	YES (+0.0002)
<i>de/wortliste</i>	0.9452	47397	0.9775	cshyphen	47544	NO (−0.0323)
<i>el/wiktionary</i>	0.9880	3524	0.9875	cshyphen	3602	YES (+0.0005)
<i>es/wiktionary</i>	0.9998	1902	0.9998	cshyphen	2023	NO (−0.0000)
<i>is/hyphenation-is</i>	0.9825	27324	0.9839	cshyphen	32284	NO (−0.0014)
<i>it/wiktionary</i>	0.9971	4158	0.9970	cshyphen	3593	YES (+0.0001)
<i>ms/wiktionary</i>	0.9302	711	0.9296	cshyphen	796	YES (+0.0006)
<i>nl/wiktionary</i>	0.9972	30097	0.9972	cshyphen	30859	YES (+0.0000)
<i>pl/wiktionary</i>	0.9915	15062	0.9913	cshyphen	8564	YES (+0.0002)
<i>pt/wiktionary</i>	0.9971	1925	0.9970	wortliste	1422	YES (+0.0001)
<i>ru/wiktionary</i>	0.9252	101559	0.9270	cshyphen	53820	NO (−0.0018)
<i>th/orchid</i>	0.9681	38955	0.9622	cshyphen	28795	YES (+0.0059)
<i>tr/wiktionary</i>	0.9915	2064	0.9920	cshyphen	1406	NO (−0.0005)
<i>uk/wiktionary</i>	0.9615	4514	0.9602	cshyphen	4741	YES (+0.0013)

References

- Franklin M. Liang and Peter Breitenlohner. 2020. **PATtern GENERation program for the \TeX 82 hyphenator.** 508–509
- Gosse Bouma. 2003. **Finite State Methods for Hyphenation.** *Natural Language Engineering*, 9(1):5–20. 510
- Jan Franěk. 2025. **Transformer-Based Word Hyphenation.** Master’s thesis, Masaryk University, Brno, Faculty of Informatics. 511–513
- Ondřej Metelka and Petr Sojka. 2025. **Hyph-bench: Benchmark Dataset of Hyphenated Words for Generating Hyphenation Patterns.** In *Proceedings of Recent Advances in Slavonic Natural Language Processing, RASLAN 2025*, pages 173–185, Brno, Czech Republic. Tribun EU. 514–516
- Lester Ingber. 1989. **Very fast simulated re-annealing.** *Mathematical and Computer Modelling*, 12(8):967–973. 517
- Carl Edward Rasmussen and Christopher K. I. Williams. 2005. **Gaussian Processes for Machine Learning (Adaptive Computation and Machine Learning).** The MIT Press. 518–519
- Scott Kirkpatrick, C. Daniel Gelatt Jr, and Mario P. Vecchi. 1983. **Optimization by simulated annealing.** *Science*, 220(4598):671–680. 520
- Donald E. Knuth. 1986. **The \TeX book.** Addison-Wesley, Reading, Massachusetts. See Appendix H for the hyphenation algorithm. 521–522
- Cornelis J. Van Rijsbergen. 1979. **Information Retrieval.** Butterworths. 523
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. **Neural Machine Translation of Rare Words with Subword Units.** In *Proceedings of the 54th Annual Meeting of the ACL (Volume 1: Long Papers)*, pages 1715–1725. ACL. 524–526
- Yan Shao, Christian Hardmeier, and Joakim Nivre. 2018. **Universal Word Segmentation: Implementation and Interpretation.** *Transactions of the Association for Computational Linguistics*, 6:421–435. 527–531
- Werner Lemberg. 2025. **A database of German words with hyphenation information.** Wordlist version with last change of Sep 20th, 2025. 532
- Jasper Snoek, Hugo Larochelle, and Ryan P. Adams. 2012. **Practical Bayesian Optimization of Machine Learning Algorithms.** In *Proceedings of the 25th International Conference on Neural Information Pro-* 533–534
- Franklin M. Liang. 1983. **Word Hy-phen-a-tion by Com-put-er.** Ph.D. thesis, Dept. of Computer Science, Stanford University. 535

536 *cessing Systems (NIPS'12)*, volume 25, pages 2951–
537 2959. Curran Associates Inc.

538 Ondřej Sojka. 2025. Transfer Learning of Slavic Syl-
539 labification for Hyphenation Patterns. Bachelor The-
540 sis, Masaryk University, Brno, Faculty of Informat-
541 ics (advisor: Pavel Šmerk), [https://is.muni.cz/
542 th/s17d6/?lang=en](https://is.muni.cz/th/s17d6/?lang=en).

543 Ondřej Sojka and Petr Sojka. 2023. Towards Perfection
544 of Machine Learning of Competing Patterns: The
545 Use Case of Czechoslovak Patterns Development. In
546 *Proceedings of the 17th Workshop on Recent Ad-
547 vances in Slavonic Natural Language Processing,
548 RASLAN 2023, Kouty nad Desnou, Czech Republic,
549 December 8–10, 2023*, pages 113–120, Brno. Tribun
550 EU.

551 Petr Sojka. 1995. Notes on Compound Word Hyphen-
552 ation in \TeX . *TUGboat*, 16(3):290–297.

553 Petr Sojka and Ondřej Sojka. 2019a. The Unreason-
554 able Effectiveness of Pattern Generation. *TUGboat*,
555 40(2):187–193.

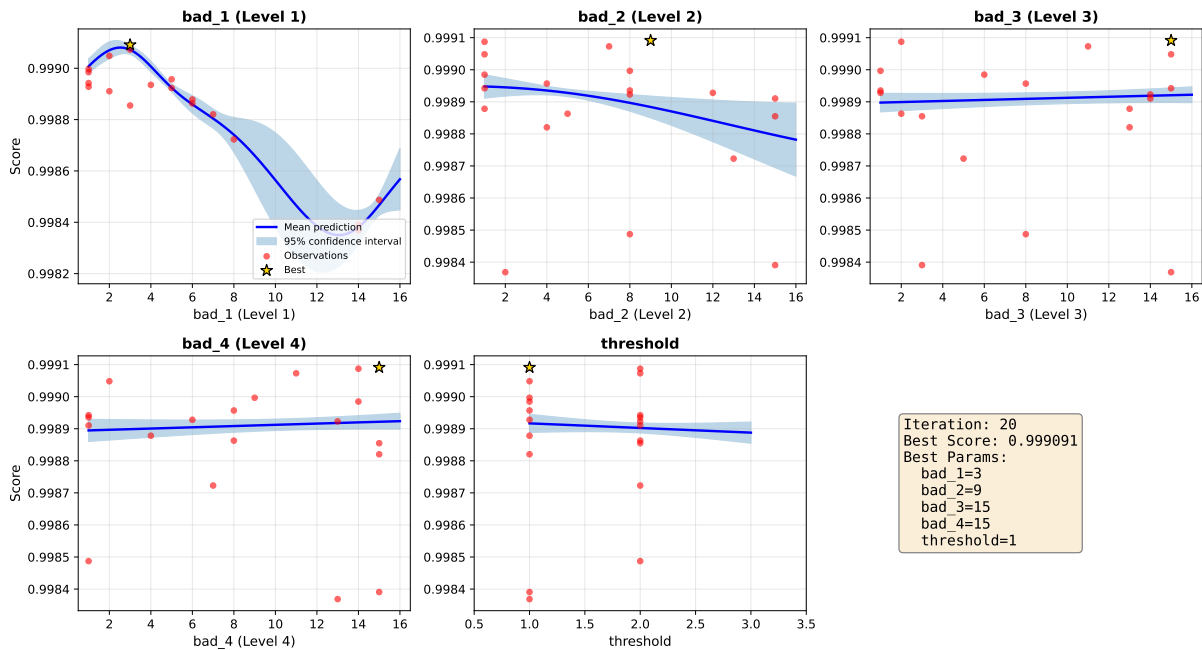
556 Petr Sojka and Ondřej Sojka. 2019b. Towards Uni-
557 versal Hyphenation Patterns. In *Proceedings of
558 Recent Advances in Slavonic Natural Language
559 Processing—RASLAN 2019*, pages 63–68, Karlova
560 Studánka, Czech Republic. Tribun EU. [https://
561 is.muni.cz/publication/1585259/?lang=en](https://is.muni.cz/publication/1585259/?lang=en).

562 Petr Sojka and Ondřej Sojka. 2021. New Czechoslovak
563 Hyphenation Patterns, Word Lists, and Workflow.
564 *TUGboat*, 42(2).

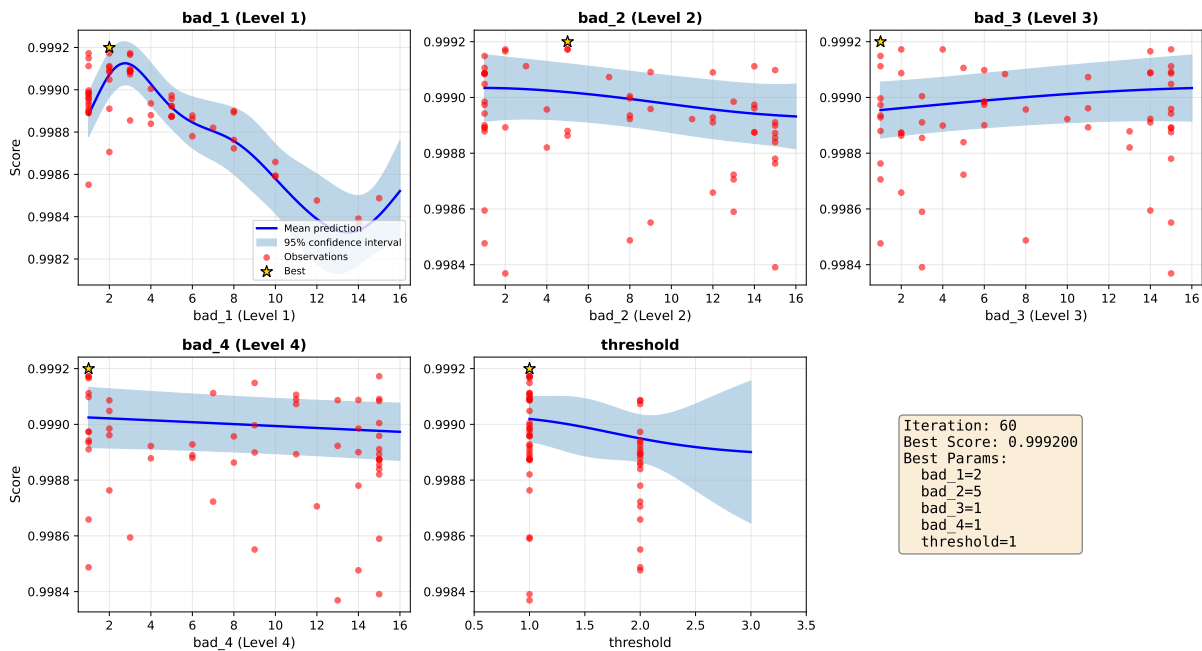
565 **A Visualizations of Gaussian Processes**

566 To better understand the parameters and results of
567 Gaussian Processes, we provide visualizations in
568 figures 2, 3, 4, and 5.

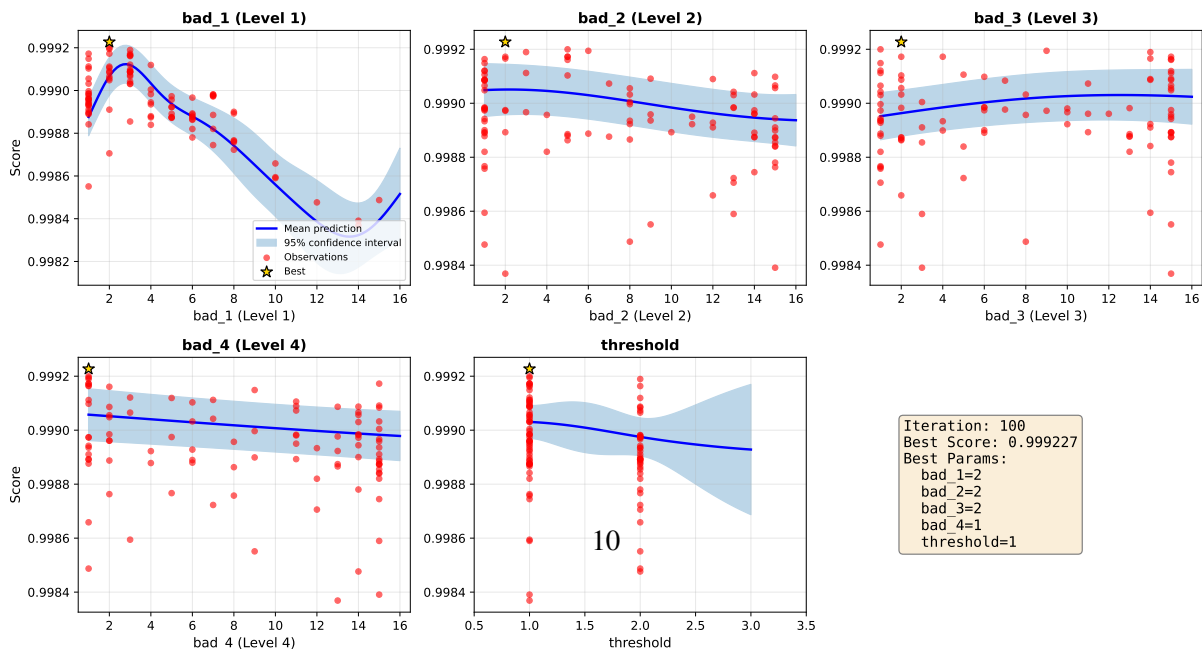
GP Regression at Iteration 20



GP Regression at Iteration 60



GP Regression at Iteration 100



GP Regression Evolution Across Iterations

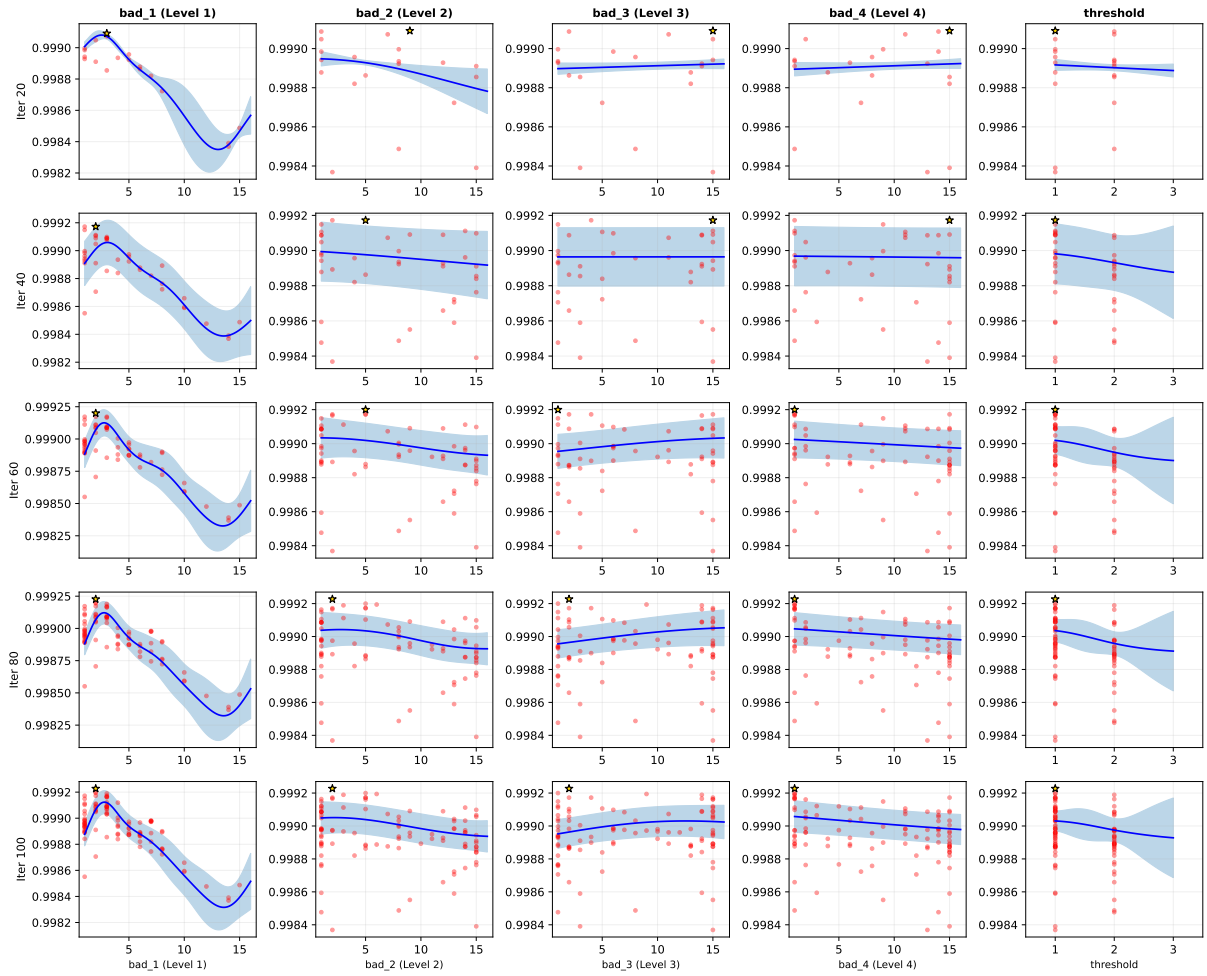


Figure 3: GP iterations 20, 60 and 100 for *csk* dataset.

GP Surface: Pairwise Parameter Interactions (Iteration 100)
Best: 0.999227 at [2, 2, 2, 1]

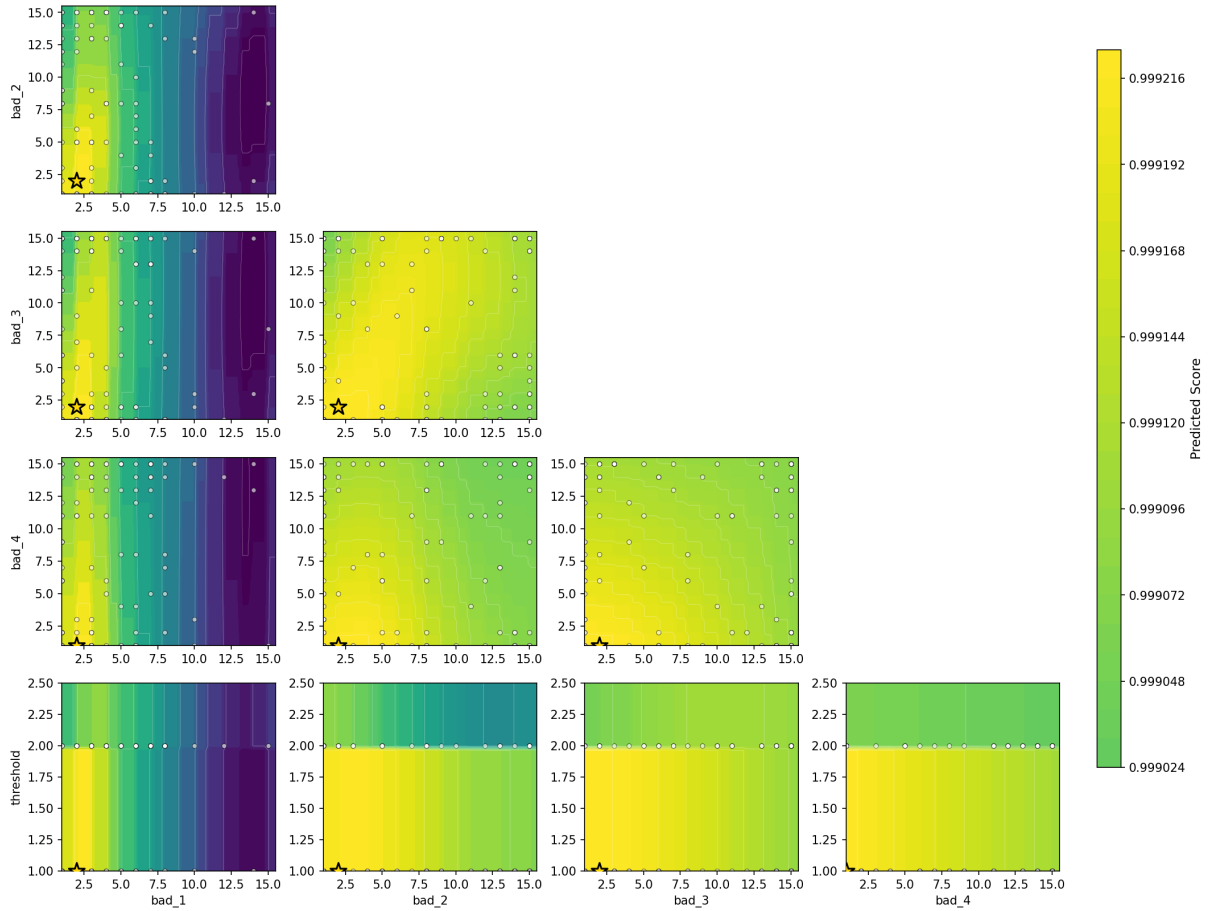


Figure 4: GP pairwise parameter interactions at iteration 100 for *cssk* dataset. “Stars” represent the optimum found.

GP Surface: Mean Prediction & Uncertainty (Iteration 100)
Best Score: 0.999227

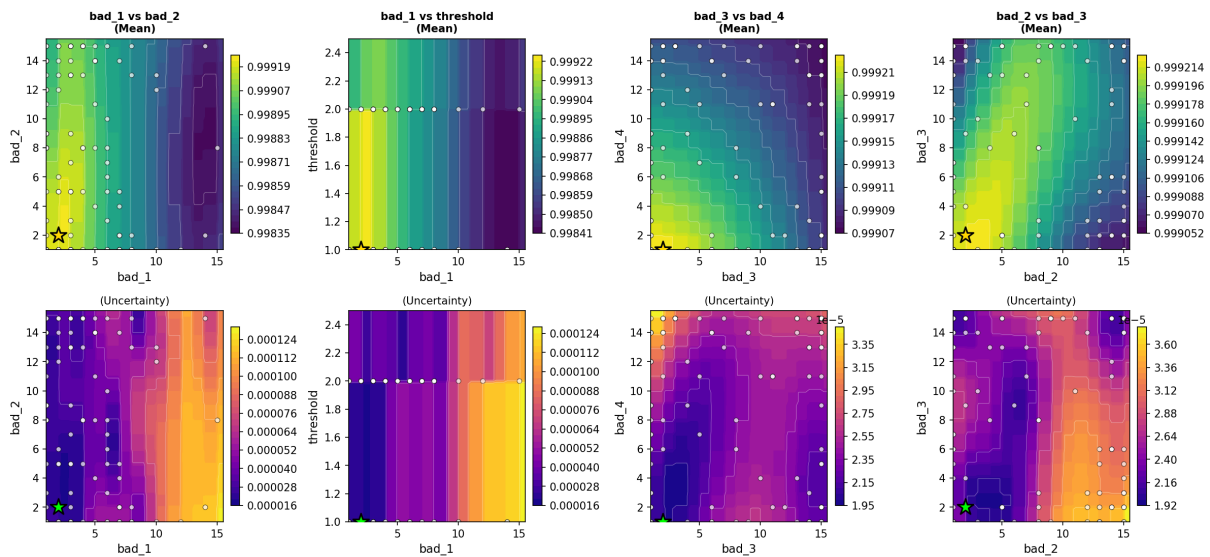


Figure 5: GP mean prediction and uncertainty at iteration 100 for *cssk* dataset.