# Space and Time-Efficient Data Structures for Massive Datasets

Giulio Ermanno Pibiri giulio.pibiri@di.unipi.it

> Supervisor Rossano Venturini

Computer Science Department University of Pisa

10/10/2017

### **High Level Thesis**

#### Data Structures + Data Compression - Faster Algorithms

#### Design **space-efficient** *ad-hoc* data structures, both from a theoretical *and* practical perspective, that support **fast data extraction**.

Data compression & Fast Retrieval together.

1. Clustered Elias-Fano Indexes

2. Dynamic Elias-Fano Representation

3. Efficient Data Structures for Massive N-Gram Datasets

#### 1. Clustered Elias-Fano Indexes

Journal paper

Giulio Ermanno Pibiri and Rossano Venturini ACM Transactions on Information Systems (TOIS), 2017

2. Dynamic Elias-Fano Representation

#### 3. Efficient Data Structures for Massive N-Gram Datasets



#### 2. Dynamic Elias-Fano Representation

Clustered Elias-Fano Indexes

1.

**Conference** paper

Giulio Ermanno Pibiri and Rossano Venturini Annual Symposium on Combinatorial Pattern Matching (CPM), 2017

#### 3. Efficient Data Structures for Massive N-Gram Datasets

#### 1. Clustered Elias-Fano Indexes

Giulio Ermanno Pibiri and Rossano Venturini ACM Transactions on Information Systems (TOIS), 2017

#### 2. Dynamic Elias-Fano Representation

Giulio Ermanno Pibiri and Rossano Venturini Annual Symposium on Combinatorial Pattern Matching (CPM), 2017

#### 3. Efficient Data Structures for Massive N-Gram Datasets

Giulio Ermanno Pibiri and Rossano Venturini ACM Conference on Research and Development in Information Retrieval (SIGIR), 2017

Conference paper



Journal paper

Giulio Ermanno Pibiri and Rossano Venturini ACM Transactions on Information Systems (TOIS), 2017

#### 2. **Dynamic Elias-Fano Representation**

Clustered Elias-Fano Indexes

1

Giulio Ermanno Pibiri and Rossano Venturini Annual Symposium on Combinatorial Pattern Matching (CPM), 2017

#### 3. Efficient Data Structures for Massive N-Gram Datasets

Giulio Ermanno Pibiri and Rossano Venturini ACM Conference on Research and Development in Information Retrieval (SIGIR), 2017

**EVERYTHING** that I do (papers, slides *and* **code**) is fully accessible at my page: http://pages.di.unipi.it/pibiri/

Journal paper

**Conference** paper









Every encoder represents each sequence individually.

No exploitation of redundancy.



Every encoder represents each sequence individually.

No exploitation of redundancy.



Idea: encode **clusters** of posting lists.

#### cluster of posting lists

:	 		
		_	
	_		
;	 		

cluster of posting lists

	,
I	
i la	
1	
·	•
· · · · · · · · · · · · · · · · · · ·	
	I
!	
	I
·	
	I
	i
	i
Ī	
	i





cluster of posting lists










- 1. Build the clusters.
- 2. Synthesise the reference list.

**NP-hard problem** already for a simplified formulation.



Figure 2: Bits per posting of Gov2 and ClueWeb09 by varying the reference size.



Figure 3: Timings for AND queries by varying the reference size on Gov2 and ClueWeb09, using the query set TREC 06.

	MIN	MID	MAX		1	MIN	MID	MAX
PEF	<b>2.94</b> (+5.60%)	<b>2.94</b> (+7.91%)	<b>2.94</b> (+10.95%)	PEF	4.80	(+2.13%)	4.80 (+3.98%)	<b>4.80</b> (+6.25%)
CPEF	2.78	2.72	2.65	CPEF	4.70		4.62	4.52
BIC	<b>2.80</b> (+0.53%)	2.80 (+2.74%)	<b>2.80</b> (+5.63%)	BIC	4.27	(-9.22%)	<b>4.27</b> (-7.58%)	<b>4.27</b> (-5.56%)
	(a)	) Gov2				(b) CI	ueWeb09	

Table 2: Bits per posting in selected trade-off points.

			MIN	Ν	MID	MAX				ИIN	MID	MAX
35	PEF	14.6	(17.5%)	14.6	(29.0%)	<b>14.6</b> (-49.7%	02	PEF	3.7	(30.4%)	<b>3.7</b> (-37.5%)	<b>3.7</b> (-52.1%)
SEC (	CPEF	17.7		20.6		29.1	SEC (	CPEF	5.3		5.9	7.8
F	BIC	41.1	(+131.9%)	41.1	(+99.5%)	<b>41.1</b> (+41.3%	Ē	BIC	10.5	(+96.2%)	10.5 (+76.2%)	10.5 (+35.0%)
90	PEF	17.7	(16.6%)	17.7	(29.1%)	17.7 (-50.3%	9	PEF	6.1	(27.4%)	<b>6.1</b> (-35.2%)	<b>6.1</b> (-49.1%
SEC (	CPEF	21.2		25.0		35.6	SEC (	CPEF	8.3		9.3	11.9
Ħ	BIC	55.1	(+159.7%)	55.1	(+120.8%)	<b>55.1</b> (+54.7%		BIC	18.5	(+122.6%)	$18.5 \ (+98.6\%)$	18.5 (+56.0%)
			(a) Cl	ueWel	b09					(b)	Gov2	

Table 3: Timings in milliseconds for AND queries on ClueWeb09 and Gov2, using query sets TREC 05 and TREC 05. In parentheses we show the relative percentage against CPEF.



Figure 2: Bits per posting of Gov2 and ClueWeb09 by varying the reference size.



Figure 3: Timings for AND queries by varying the reference size on Gov2 and ClueWeb09, using the query set TREC 06.

	MIN	MID	МАХ		MIN	MID	МАХ
PEF	<b>2.94</b> (+5.60%)	<b>2.94</b> (+7.919	) <b>2.94</b> (+10.95%)	PEF	<b>4.80</b> (+2.13%)	<b>4.80</b> (+3.98	6) <b>4.80</b> (+6.25%)
CPEF	2.78	2.72	2.65	CPEF	4.70	4.62	4.52
BIC	<b>2.80</b> (+0.53%)	2.80 (+2.74	) <b>2.80</b> (+5.63%)	BIC	<b>4.27</b> (-9.22%)	<b>4.27</b> (-7.58	b) <b>4.27</b> (-5.56%)
	(a)	) Gov2			(b) CI	ueWeb09	

Table 2: Bits per posting in selected trade-off points.

		1	MIN	N	٨ID	MAX			1	MIN	MID	MAX
35	PEF	14.6	(17.5%)	14.6	(29.0%)	14.6 (-49.7%)	5	PEF	3.7	(30.4%)	<b>3.7</b> (-37.5%)	<b>3.7</b> (-52.1%)
SEC (	CPEF	17.7		20.6		29.1	SEC (	CPEF	5.3		5.9	7.8
Ŧ	BIC	41.1	(+131.9%)	41.1	(+99.5%)	<b>41.1</b> (+41.3%)	Ĩ	BIC	10.5	(+96.2%)	10.5 (+76.2%)	10.5 (+35.0%)
90	PEF	17.7	(16.6%)	17.7	(29.1%)	17.7 (-50.3%)	90	PEF	6.1	(27.4%)	<b>6.1</b> (-35.2%)	<b>6.1</b> (-49.1%)
SEC (	CPEF	21.2		25.0		35.6	SEC (	CPEF	8.3		9.3	11.9
F	BIC	55.1	(+159.7%)	55.1	(+120.8%)	55.1 (+54.7%)	F	BIC	18.5	(+122.6%)	$18.5 \ (+98.6\%)$	18.5 (+56.0%)
			(a) Cli	JeWel	o09					(b)	Gov2	

Table 3: Timings in milliseconds for AND queries on ClueWeb09 and Gov2, using query sets TREC 05 and TREC 05. In parentheses we show the relative percentage against CPEF.



Figure 2: Bits per posting of Gov2 and ClueWeb09 by varying the reference size.



Figure 3: Timings for AND queries by varying the reference size on Gov2 and ClueWeb09, using the query set TREC 06.

	MIN	MID	МАХ		MIN	MID	МАХ
PEF	<b>2.94</b> (+5.60%)	<b>2.94</b> (+7.919	) <b>2.94</b> (+10.95%)	PEF	<b>4.80</b> (+2.13%)	4.80 (+3.98	6) <b>4.80</b> (+6.25%)
CPEF	2.78	2.72	2.65	CPEF	4.70	4.62	4.52
BIC	<b>2.80</b> (+0.53%)	2.80 (+2.74	) <b>2.80</b> (+5.63%)	BIC	<b>4.27</b> (-9.22%)	<b>4.27</b> (-7.58	6) <b>4.27</b> (-5.56%)
	(a)	) Gov2			(b) CI	ueWeb09	

Table 2: Bits per posting in selected trade-off points.

### Always better than PEF (by up to 11%) and better than BIC (by up to 6.25%)

_												
			MIN	N	MID	MAX				MIN	MID	MAX
35	PEF	14.6	(17.5%)	14.6	(29.0%)	<b>14.6</b> (-49.7%	) 6	PEF	3.7	(30.4%)	<b>3.7</b> (-37.5%)	<b>3.7</b> (-52.1%)
SEC (	CPEF	17.7		20.6		29.1	SEC (	CPEF	5.3		5.9	7.8
Ē	BIC	41.1	(+131.9%)	41.1	(+99.5%)	<b>41.1</b> (+41.3%		BIC	10.5	(+96.2%)	10.5 (+76.2%)	10.5 (+35.0%)
90	PEF	17.7	(16.6%)	17.7	(29.1%)	17.7 (-50.3%	) 9	PEF	6.1	(27.4%)	<b>6.1</b> (-35.2%)	<b>6.1</b> (-49.1%)
REC	CPEF	21.2		25.0		35.6	REC	CPEF	8.3		9.3	11.9
F	BIC	55.1	(+159.7%)	55.1	(+120.8%)	55.1 (+54.7%		BIC	18.5	(+122.6%)	$18.5 \ (+98.6\%)$	18.5 (+56.0%)
			(a) Clu	JeWel	09					(b)	Gov2	

Table 3: Timings in milliseconds for AND queries on ClueWeb09 and Gov2, using query sets TREC 05 and TREC 05. In parentheses we show the relative percentage against CPEF.



Figure 2: Bits per posting of Gov2 and ClueWeb09 by varying the reference size.



Figure 3: Timings for AND queries by varying the reference size on Gov2 and ClueWeb09, using the query set TREC 06.

	MIN	MID	MAX		MIN	MID	MAX
PEF	<b>2.94</b> (+5.60%)	<b>2.94</b> (+7.919	) <b>2.94</b> (+10.95%)	PEF	<b>4.80</b> (+2.13%)	<b>4.80</b> (+3.98	6) <b>4.80</b> (+6.25%)
CPEF	2.78	2.72	2.65	CPEF	4.70	4.62	4.52
BIC	<b>2.80</b> (+0.53%)	2.80 (+2.749	) <b>2.80</b> (+5.63%)	BIC	<b>4.27</b> (-9.22%)	<b>4.27</b> (-7.58	6) <b>4.27</b> (-5.56%)
	(a)	) Gov2			(b) CI	ueWeb09	

Table 2: Bits per posting in selected trade-off points.

## Always better than PEF (by up to 11%) and better than BIC (by up to 6.25%)

		MIN	MID	MAX			N	ліN	MID	MAX
05	PEF	<b>14.6</b> (-17.5%)	4.6 (-29.0%)	14.6 (-49.7%)	05	PEF	3.7	(30.4%)	<b>3.7</b> (-37.5%)	<b>3.7</b> (-52.1%)
SEC	CPEF	17.7	20.6	29.1	Ê	CPEF	5.3		5.9	7.8
Ţ	BIC	<b>41.1</b> (+131.9%)	1.1 (+99.5%)	<b>41.1</b> (+41.3%)	Ĩ	віс	10.5	(+96.2%)	10.5 (+76.2%)	10.5 (+35.0%)
90	PEF	<b>17.7</b> (-16.6%)	<b>.7.7</b> (-29.1%)	17.7 (-50.3%)	90	PEF	6.1	(27.4%)	<b>6.1</b> (-35.2%)	<b>6.1</b> (-49.1%)
REC	CPEF	21.2	25.0	35.6	ZEC	CPEF	8.3		9.3	11.9
F	BIC	55.1 (+159.7%)	<b>5.1</b> (+120.8%)	<b>55.1</b> (+54.7%)	F	віс	18.5	(+122.6%)	18.5 (+98.6%)	$18.5 \ (+56.0\%)$
_		(a) Clu	eWeb09					(b)	Gov2	

Table 3: Timings in milliseconds for AND queries on ClueWeb09 and Gov2, using query sets TREC 05 and TREC 05. In parentheses we show the relative percentage against CPEF.



Figure 2: Bits per posting of Gov2 and ClueWeb09 by varying the reference size.



Figure 3: Timings for AND queries by varying the reference size on Gov2 and ClueWeb09, using the query set TREC 06.

	MIN	MID	МАХ		MIN	MID	MAX
PEF	<b>2.94</b> (+5.60%)	<b>2.94</b> (+7.919	) <b>2.94</b> (+10.95%)	PEF	<b>4.80</b> (+2.13%)	4.80 (+3.98	6) <b>4.80</b> (+6.25%)
CPEF	2.78	2.72	2.65	CPEF	4.70	4.62	4.52
BIC	<b>2.80</b> (+0.53%)	2.80 (+2.74	) <b>2.80</b> (+5.63%)	BIC	<b>4.27</b> (-9.22%)	<b>4.27</b> (-7.58	6) <b>4.27</b> (-5.56%)
	(a)	) Gov2			(b) CI	ueWeb09	

Table 2: Bits per posting in selected trade-off points.

Always better than PEF (by up to 11%) and better than BIC (by up to 6.25%)

		MI	N	Μ	IID	MA	X			N	/IN	MID	MAX
es PE	F 1	4.6 (-	-17.5%)	. <b>4.</b> 6	(29.0%)	<b>14.6</b> (	49.7%)	05	PEF	3.7	(30.4%)	<b>3.7</b> (-37.5%)	<b>3.7</b> (-52.1%)
Щ CP	PEF 1	7.7		20.6		29.1		SEC	CPEF	5.3		5.9	7.8
Ёвк	c 4	<b>41.1</b> (+131.9%)		<b>1.1</b> (+99.5%)		<b>41.1</b> (+41.3%)		Ŧ	віс	10.5	(+96.2%)	10.5 (+76.2%)	10.5 (+35.0%)
8 PE	F 1	17.7 (- 21.2	(-16.6%) . <b>7.7</b> . <b>25.0</b>	.7.7	(—29.1%)	$\begin{array}{c} 17.7 (-50.3\%) \\ 35.6 \end{array}$	50.3%)	3EC 06	PEF	6.1	(27.4%)	<b>6.1</b> (-35.2%)	<b>6.1</b> (-49.1%)
ы СР	PEF 2			25.0					CPEF	8.3		9.3	11.9
F BIC	c <b>5</b>	5.1 (+1	159.7%)	<b>5.1</b> (	+120.8%)	55.1 (+4	54.7%)	F	віс	18.5	(+122.6%)	18.5 (+98.6%)	18.5 (+56.0%)
(a) ClueWeb09						(b) Gov2							

Much faster than BIC (103% on average) Slightly slower than PEF (20% on average)

### (Integer) Dynamic Ordered Sets

A dynamic ordered set S is a data structure representing n

keys and supporting the following operations:

- Insert(x) inserts x in S
- Delete(x) deletes x from S
- Search(x) checks whether x belongs to S
- Minimum() returns the minimum element of S
- Maximum() returns the maximum element of S
- Predecessor(x) returns  $max{y \in S : y < x}$
- Successor(x) returns  $\min\{y \in S : y \ge x\}$

### (Integer) Dynamic Ordered Sets

A dynamic ordered set S is a data structure representing n

keys and supporting the following operations:

- Insert(x) inserts x in S
- Delete(x) deletes x from S
- Search(x) checks whether x belongs to S
- Minimum() returns the minimum element of S
- Maximum() returns the maximum element of S
- Predecessor(x) returns  $\max{y \in S : y < x}$
- Successor(x) returns  $\min\{y \in S : y \ge x\}$

In the **comparison model** this is solved optimally by any self-balancing tree data structure in  $O(\log n)$  time and O(n) space.

More efficient solutions there exist if the considered keys are **integers** drawn from a bounded universe of size *u*.

### (Integer) Dynamic Ordered Sets

A dynamic ordered set *S* is a data structure representing *n* 

keys and supporting the following operations:

- Insert(x) inserts x in S
- Delete(x) deletes x from S
- Search(x) checks whether x belongs to S
- Minimum() returns the minimum element of S
- Maximum() returns the maximum element of S
- Predecessor(x) returns  $\max\{y \in S : y < x\}$
- Successor(x) returns  $\min\{y \in S : y \ge x\}$

In the **comparison model** this is solved optimally by any self-balancing tree data structure in  $O(\log n)$  time and O(n) space.

More efficient solutions there exist if the considered keys are **integers** drawn from a bounded universe of size *u*.

#### Challenge How to **optimally** solve the **integer** dynamic

ordered set problem in **compressed space**?

### **Motivation**

#### **Integer Data Structures**

- van Emde Boas Trees
- X/Y-Fast Tries
- Fusion Trees
- Exponential Search Trees

- + time
- space
- + dynamic

#### **Elias-Fano Encoding**

 $EF(S(n,u)) = n \log(u/n) + 2n$  bits to

encode an ordered integer sequence *S* 

- O(1) Access
- $O(1 + \log(u/n))$  **Predecessor** 
  - + time
  - + space
  - static

### **Motivation**

#### **Integer Data Structures**

- van Emde Boas Trees
- X/Y-Fast Tries
- Fusion Trees
- Exponential Search Trees

- + time
- space
- + dynamic

#### **Elias-Fano Encoding**

 $EF(S(n,u)) = n \log(u/n) + 2n$  bits to

encode an ordered integer sequence *S* 

- O(1) Access
- $O(1 + \log(u/n))$  **Predecessor** 
  - + time
  - + space
  - static

### **Motivation**



I

For $u = n^{\gamma}$ , $\gamma = \Theta(1)$ :	
• $EF(S(n,u)) + O(n)$ bits	Result 1
• O(1) Access	
<ul> <li>O(min{1+log(u/n), loglog n}) Predecessor</li> </ul>	
• $EF(S(n,u)) + O(n)$ bits	
• O(1) Access	Result 2
<ul> <li>O(1) Append (amortized)</li> </ul>	
<ul> <li>O(min{1+log(u/n), loglog n}) Predecessor</li> </ul>	
	1
• $EF(S(n,u)) + O(n)$ bits	
• $O(\log n / \log \log n)$ Access	Rocult 3
<ul> <li>O(log n / loglog n) Insert/Delete (amortized)</li> </ul>	
<ul> <li>O(min{1+log(u/n), loglog n}) Predecessor</li> </ul>	

Fc • •	or $u = n^{\gamma}, \gamma = \Theta(1)$ : EF( $S(n,u)$ ) + O( $n$ ) bits O(1) Access O(min{1+log( $u/n$ ), loglog $n$ }) Predecessor	Result 1
•	$ EF(S(n,u)) + O(n) bits  O(1) Access  O(1) Append (amortized)  O(min{1+log(u/n), loglog n}) Predecessor$	Result 2
•	$EF(S(n,u)) + o(n) bits$ $O(\log n / \log\log n) Access$ $O(\log n / \log\log n) Insert/Delete (amortized)$ $O(\min\{1+\log(u/n), \log\log n\}) Predecessor$	Result 3

For $u = n^{\gamma}$ , $\gamma = \Theta(1)$ : • EF( $S(n,u)$ ) + $O(n)$ bits • O(1) Access • O(min{1+log( $u/n$ ), loglog $n$ }) Predecessor	Result 1
<ul> <li>EF(S(n,u)) + o(n) bits</li> <li>O(1) Access</li> <li>O(1) Append (amortized)</li> <li>O(min{1+log(u/n), loglog n}) Predecessor</li> </ul>	Result 2
<ul> <li>EF(S(n,u)) + o(n) bits</li> <li>O(log n / loglog n) Access</li> <li>O(log n / loglog n) Insert/Delete (amortized)</li> <li>O(min{1+log(u/n), loglog n}) Predecessor</li> </ul>	Result 3




























Strings of *N* words. *N* typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.

Strings of N words. N typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.



Strings of N words. N typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.



# Google Books

 $\approx$  6% of the books ever published

Strings of N words. N typically ranges from 1 to 5.

Extracted from text using a *sliding window* approach.



# Google Books

 $\approx$  6% of the books ever published

N	number of grams
1	24,359,473
2	667,284,771
3	7,397,041,901
4	1,644,807,896
5	1,415,355,596

More than 11 billion grams.











- Millions of unigrams.
- Height 5: longer contexts.
- The number of siblings has a **funnel-**shaped distribution.

- Millions of unigrams.
- Height 5: longer contexts.
- The number of siblings has a **funnel-**shaped distribution.



- Millions of unigrams.
- Height 5: longer contexts.
- The number of siblings has a **funnel-**shaped distribution.



u/n by varying context-length k

	$\boldsymbol{k}$	3-grams	4-grams	5-grams
Europarl	0 1 2	$2404 \\ 213 \ (\times 11.28) \\ 2404$	2782 480 (×5.79) 48 (×57.95)	2920 646 (×4.52) 101 (×28.91)
YahooV2	0 1 2	7350 753 (×9.76) 7350	$\begin{array}{c} 7197 \\ 1461  (\times 4.93) \\ 104 \ (\times 69.20) \end{array}$	$\begin{array}{c} 7417 \\ 1963  (\times 3.78) \\ 249 \ (\times 29.79) \end{array}$
GoogleV2	0 1 2	4050 1025 (×3.95) 4050	$\begin{array}{c} 6631 \\ 2192  (\times 3.03) \\ 221 \ (\times 30.00) \end{array}$	6793 2772 (×2.45) 503 (×13.50)

- Millions of unigrams.
- Height 5: longer contexts.
- The number of siblings has a **funnel-**shaped distribution.



u/n by varying context-length k

	$\boldsymbol{k}$	3-grams	4-grams	5-grams
arl	0	2404	2782	2920
rop	1	<b>213</b> (×11.28)	480 (×5.79)	646 (×4.52)
Eu	2	2404	$48_{(\times 57.95)}$	$101 \ (\times 28.91)$
<b>V2</b>	0	7350	7197	7417
oot	1	753 (×9.76)	1461 (×4.93)	1963 (×3.78)
Yał	2	7350	$104_{( imes 69.20)}$	$249 \ (\times 29.79)$
22	0	4050	6631	6793
ogle	1	1025 (×3.95)	2192 (×3.03)	2772 (×2.45)
go	2	4050	$221_{( imes 30.00)}$	$503 (\times 13.50)$

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665752080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
5	43160518	295701337	1413870914
Total	101428257	828 223 677	11131242087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304 579	3475482	24 357 349
2	5192260	53844927	665752080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
5	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665 752 080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
5	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

	Eur	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$	
EF PEF	1.97 1.87 (-4.99%)	1.28 1.35 (+5.93%)	2.17 1.91 (-12.03%)	1.60 1.73 (+8.00%)	$\begin{array}{c} \textbf{2.13} \\ \textbf{1.52} \ (-28.60\%) \end{array}$	2.09 1.91 (-8.79%)	
$\begin{array}{c} \textbf{T-BASED} \\ \textbf{APPING} \\ \textbf{k} = 1 \\ \textbf{APL} \\ APL$	$\frac{1.67}{1.53} \left( ^{-15.30\%} \right) \\ \frac{1.53}{(-22.36\%)}$	$\begin{array}{c} \textbf{1.58} \\ \textbf{(+23.86\%)} \\ \textbf{1.61} \\ \textbf{(+25.89\%)} \end{array}$	$\frac{1.89}{1.63} \left( ^{-12.92\%} \right)$	$\begin{array}{c} \textbf{2.05} \scriptstyle{(+28.07\%)} \\ \textbf{2.16} \scriptstyle{(+35.22\%)} \end{array}$	$\begin{array}{c} \textbf{1.91} \ (-10.24\%) \\ \textbf{1.31} \ (-38.71\%) \end{array}$	3.03 (+44.61%) 2.30 (+9.88%)	
$\begin{bmatrix} CONTEX \\ ID REM \\ k = 2 \\ H = 2 \\ $	$\frac{1.46}{1.28} \left( ^{-25.62\%} \right) \\ \left( ^{-34.87\%} \right)$	$\frac{1.60}{1.64}_{(+28.12\%)}$	$\frac{1.68}{1.38} \left( ^{-22.32\%} \right) \\ \frac{1.38}{(-36.15\%)}$	$\begin{array}{c} \textbf{2.08} (+30.23\%) \\ \textbf{2.15} (+34.81\%) \end{array}$	_	_	

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665752080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
<b>5</b>	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

	Eur	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$	
EF PEF	1.97 1.87 (-4.99%)	1.28 1.35 (+5.93%)	2.17 1.91 (-12.03%)	1.60 1.73 (+8.00%)	$\begin{array}{c} \textbf{2.13} \\ \textbf{1.52} \ (-28.60\%) \end{array}$	2.09 1.91 (-8.79%)	
$\begin{array}{c} \text{T-BASED} \\ \text{APPING} \\ k = 1 \\ \text{Herbits} \\ He$	$\frac{1.67}{1.53}_{(-22.36\%)}$	$\begin{array}{c} \textbf{1.58} \\ \textbf{(+23.86\%)} \\ \textbf{1.61} \\ \textbf{(+25.89\%)} \end{array}$	$\frac{1.89}{1.63}_{(-24.91\%)}$	$\begin{array}{c} \textbf{2.05} (+28.07\%) \\ \textbf{2.16} (+35.22\%) \end{array}$	$\frac{1.91}{1.31}_{(-38.71\%)}^{(-10.24\%)}$	3.03 (+44.61%) 2.30 (+9.88%)	
	$\frac{1.46}{1.28}_{(-34.87\%)}^{(-25.62\%)}$	$\frac{1.60}{1.64}_{(+28.12\%)}$	$\frac{1.68}{1.38}_{(-36.15\%)}^{(-22.32\%)}$	$\begin{array}{c} \textbf{2.08} (+30.23\%) \\ \textbf{2.15} (+34.81\%) \end{array}$	_	_	

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665 752 080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
<b>5</b>	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

**C++** implementation gcc 5.4.1 with the highest optimization setting

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF PEF	1.97 1.87 (-4.99%)	1.28 1.35 (+5.93%)	2.17 1.91 (-12.03%)	1.60 1.73 (+8.00%)	$\begin{array}{c} \textbf{2.13} \\ \textbf{1.52} \ (-28.60\%) \end{array}$	2.09 1.91 (-8.79%)
$\begin{array}{c} \text{T-BASED} \\ \text{APPING} \\ k = 1 \\ \text{APPING} \\ $	$\frac{1.67}{1.53}_{(-22.36\%)}^{(-15.30\%)}$	$\begin{array}{c} \textbf{1.58} \\ \textbf{(+23.86\%)} \\ \textbf{1.61} \\ \textbf{(+25.89\%)} \end{array}$	$\frac{1.89}{1.63}_{(-24.91\%)}$	$\begin{array}{c} \textbf{2.05} \scriptstyle{(+28.07\%)} \\ \textbf{2.16} \scriptstyle{(+35.22\%)} \end{array}$	$\frac{1.91}{1.31}_{(-38.71\%)}^{(-10.24\%)}$	$\begin{array}{c} \textbf{3.03} (+44.61\%) \\ \textbf{2.30} (+9.88\%) \end{array}$
$\begin{array}{c} \text{CONTEX} \\ \text{CONTEX} \\ \text{ID REM} \\ \text{F} = 2 \\ \text{F} \\ \text{F} = 2 \\ \text{F} \\ $	$\frac{1.46}{1.28}_{(-34.87\%)}^{(-25.62\%)}$	$\frac{1.60}{1.64}_{(+28.12\%)}$	$\frac{1.68}{1.38}_{(-36.15\%)}^{(-22.32\%)}$	$\begin{array}{c} \textbf{2.08} (+30.23\%) \\ \textbf{2.15} (+34.81\%) \end{array}$	_	_

#### **Context-based ID Remapping**

reduces space by more than 36% on average ----- you will notice this!

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665752080
3	18908249	187639522	7384478110
4	33862651	287562409	1642783634
5	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

**C++** implementation gcc 5.4.1 with the highest optimization setting

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF PEF	1.97 1.87 (-4.99%)	1.28 1.35 (+5.93%)	2.17 1.91 (-12.03%)	1.60 1.73 (+8.00%)	$\begin{array}{c} \textbf{2.13} \\ \textbf{1.52} \ (-28.60\%) \end{array}$	2.09 1.91 (-8.79%)
$\begin{array}{c} \text{T-BASED} \\ \text{APPING} \\ k = 1 \\ \text{APPING} \\ $	$\frac{1.67}{1.53}_{(-22.36\%)}$	$\frac{1.58}{1.61}_{(+25.89\%)}$	$\frac{1.89}{1.63}_{(-24.91\%)}$	$2.05_{(+28.07\%)}\\2.16_{(+35.22\%)}$	$\frac{1.91}{1.31}_{(-38.71\%)}^{(-10.24\%)}$	$3.03 \atop (+44.61\%) \atop (+9.88\%)$
	$\frac{1.46}{1.28}_{(-34.87\%)}$	$\frac{1.60}{(+25.17\%)}$ 1.64 (+28.12%)	$\frac{1.68}{1.38}_{(-36.15\%)}$	$\begin{array}{c} 2.08 \\ (+30.23\%) \\ 2.15 \\ (+34.81\%) \end{array}$	_	

#### **Context-based ID Remapping**

reduces space by more than 36% on average ----- you will notice this!

N	Europarl	YahooV2	GoogleV2
	n	n	n
1	304579	3475482	24 357 349
2	5192260	53844927	665 752 080
3	18908249	187639522	7384478110
4	33862651	287 562 409	1642783634
5	43160518	295701337	1413870914
Total	101428257	828 223 677	11 131 242 087
gzip bpg	6.98	6.45	6.20

Test machine Intel Xeon E5-2630 v3, 2.4 GHz 193 GB of RAM, Linux 64 bits

**C++** implementation gcc 5.4.1 with the highest optimization setting

	Europarl		YahooV2		GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
EF PEF	1.97 1.87 (-4.99%)	1.28 1.35 (+5.93%)	2.17 1.91 (-12.03%)	1.60 1.73 (+8.00%)	$\begin{array}{c} \textbf{2.13} \\ \textbf{1.52} \ (-28.60\%) \end{array}$	2.09 1.91 (-8.79%)
$\begin{array}{c} \text{Presed} \\ \text{APPING} \\ k = 1 \\ k = 1 \\ \text{APPING} \\ \text{APPING \\ \text{APPING} \\ \text{APPING} \\ APPING$	$\frac{1.67}{1.53}_{(-22.36\%)}$	$\frac{1.58}{1.61}_{(+25.89\%)}$	$\frac{1.89}{1.63}_{(-24.91\%)}$	$2.05_{(+28.07\%)}\\2.16_{(+35.22\%)}$	$\frac{1.91}{1.31}_{(-38.71\%)}^{(-10.24\%)}$	$3.03 \atop (+44.61\%) \\ 2.30 \atop (+9.88\%)$
	$\frac{1.46}{1.28}_{(-34.87\%)}$	$\begin{array}{c} \textbf{1.60} \\ \textbf{(+25.17\%)} \\ \textbf{1.64} \\ \textbf{(+28.12\%)} \end{array}$	$\frac{1.68}{1.38}_{(-36.15\%)}$	$\begin{array}{c} \textbf{2.08} \\ \textbf{(+30.23\%)} \\ \textbf{2.15} \\ \textbf{(+34.81\%)} \end{array}$	_	_

#### **Context-based ID Remapping**

- reduces space by more than 36% on average -
- brings approximately **30%** more time

- you will notice this!
- will you notice this?

	Eur	Europarl		YahooV2		ogleV2
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie PEF-RTrie	$1.87 \\ 1.28$	$1.35 \\ 1.64$	1.91 1.38	1.73 2.15	$1.52 \\ 1.31$	1.91 2.30
BerkeleyLM C.	1.70 (-8.89%)	2.83 (+108.88%)	<b>1.69</b> (-11.41%)	3.48 (+101.84%)	1.45 (-4.87%)	4.13 (+116.57%)
BerkeleyLM H.3	(+32.90%) 6.70 (+258.81%)	(+72.70%) <b>0.97</b> (-28.46%)	(+22.04%) <b>7.82</b> (+310.38%)	(+61.70%) 1.13 (-34.35%)	(+10.83%) 9.24 (+507.79%)	(+79.76%) <b>2.18</b> (+13.95%)
BerkeleyLM H.50	(+423.40%) <b>7.96</b> (+326.03%)	(-40.85%) <b>0.97</b> (-28.49%)	(+465.36%) 9.37 (+391.32%)	(-47.41%) <b>0.96</b> (-44.27%)	(+608.07%)	(-5.42%)
Expgram	(+521.45%) <b>2.06</b> (+10.18%)	(-40.88%) <b>2.80</b> (+106.61%)	(+576.87%) <b>2.24</b> (+17.36%)	(-55.35%) 9.23 (+435.33%)		_
KenLM T.	(+60.73%) <b>2.99</b> (+60.11%)	(+70.82%) <b>1.28</b> (-5.47%)	(+61.68%) <b>3.44</b> (+80.39%)	(+328.87%) <b>1.94</b> (+12.32%)		_
Marisa	(+133.56%) <b>3.61</b> (+93.09%)	(-21.84%) <b>2.06</b> (+52.00%)	(+148.52%) <b>3.81</b> (+99.60%)	(-10.01%) <b>3.24</b> (+87.96%)	_	_
RandLM	(+181.66%) <b>1.81</b> (-3.06%)	(+25.67%) 4.39 (+224.20%)	<b>2.02</b> (+6.18%)	(+50.58%) 5.08 (+194.35%)	<b>2.60</b> (+70.73%)	9.25 (+384.54%)
	(+41.41%)	(+168.04%)	(+46.29%)	(+135.82%)	(+98.90%)	(+302.19%)

		Euro	Europarl		YahooV2		ogleV2
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie PEF-RTrie		$1.87 \\ 1.28$	$1.35 \\ 1.64$	1.91 1.38	1.73 2.15	$1.52 \\ 1.31$	1.91 2.30
BerkeleyLM	С.	1.70 (-8.89%)	2.83 (+108.88%)	1.69 (-11.41%)	3.48 (+101.84%)	1.45 (-4.87%)	4.13 (+116.57%)
BerkeleyLM	H.3	(+32.90%) 6.70 (+258.81%)	(+72.70%) <b>0.97</b> (-28.46%)	(+22.04%) <b>7.82</b> (+310.38%)	(+61.70%) 1.13 (-34.35%)	(+10.83%) 9.24 (+507.79%)	(+79.76%) <b>2.18</b> (+13.95%)
BerkeleyLM	H.50	(+423.40%) 7.96 (+326.03%)	(-40.85%) <b>0.97</b> (-28.49%)	(+465.36%) 9.37 (+391.32%)	(-47.41%) <b>0.96</b> (-44.27%)	(+608.07%)	(-5.42%)
Expgram		(+521.45%) <b>2.06</b> (+10.18%)	(-40.88%) <b>2.80</b> (+106.61%)	(+576.87%) <b>2.24</b> (+17.36%)	(-55.35%) 9.23 (+435.33%)		_
KenLM T.		<b>2.99</b> (+60.11%)	1.28 (-5.47%)	<b>3.44</b> (+80.39%)	1.94 (+12.32%)	—	_
Marisa		<b>3.61</b> (+93.09%) (+181.66%)	<b>2.06</b> (+52.00%) (+25.67%)	<b>3.81</b> (+99.60%) (+174.98%)	<b>3.24</b> (+87.96%) (+50.58%)		_
RandLM		1.81 (-3.06%)	4.39 (+224.20%)	2.02 (+6.18%)	5.08 (+194.35%)	2.60 (+70.73%)	9.25 (+384.54%)
		(+41.41%)	(+168.04%)	(+46.29%)	(+135.82%)	(+98.90%)	(+302.19%)

		Euro	Europarl		YahooV2		ogleV2
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie PEF-RTrie		$1.87 \\ 1.28$	$1.35 \\ 1.64$	1.91 1.38	1.73 2.15	$1.52 \\ 1.31$	1.91 2.30
BerkeleyLM	С.	1.70 (-8.89%)	2.83 (+108.88%)	<b>1.69</b> (-11.41%)	3.48 (+101.84%)	1.45 (-4.87%)	4.13 (+116.57%)
BerkeleyLM	H.3	(+32.90%) 6.70 (+258.81%)	(+72.70%) <b>0.97</b> (-28.46%)	(+22.04%) <b>7.82</b> (+310.38%)	(+61.70%) 1.13 (-34.35%)	(+10.83%) 9.24 (+507.79%)	(+79.76%) <b>2.18</b> (+13.95%)
BerkeleyLM	H.50	(+423.40%) 7.96 (+326.03%)	(-40.85%) <b>0.97</b> (-28.49%)	(+465.36%) 9.37 (+391.32%)	(-47.41%) <b>0.96</b> (-44.27%)	(+608.07%)	(-5.42%)
Expgram		(+521.45%) <b>2.06</b> (+10.18%)	(-40.88%) <b>2.80</b> (+106.61%)	(+576.87%) <b>2.24</b> (+17.36%)	(-55.35%) 9.23 (+435.33%)	_	_
KenLM T.		2.99 (263X <sup>6</sup> )	<b>1.28</b> (-5.47%)	3.44 (28.5%)	<b>1.94</b> (+12.32%)		_
Marisa		<b>3.61</b> (+93.09%) (+181.66%)	<b>2.06</b> (+52.00%) (+25.67%)	<b>3.81</b> (+99.60%) (+174.98%)	<b>3.24</b> (+87.96%) (+50.58%)	—	—
RandLM		1.81 (-3.06%)	4.39 (+224.20%)	2.02 (+6.18%)	5.08 (+194.35%)	<b>2.60</b> (+70.73%)	9.25 (+384.54%)
		(+41.41%)	(+168.04%)	(+46.29%)	(+135.82%)	(+98.90%)	(+302.19%)

		Europarl			YahooV2		GoogleV2	
		bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$	
PEF-Trie PEF-RTrie		$1.87 \\ 1.28$	1.35 1.64	$1.91 \\ 1.38$	$\begin{array}{c} 1.73 \\ 2.15 \end{array}$	$1.52 \\ 1.31$	1.91 2.30	
BerkeleyLM	C.	1.70 (-8.89%)	2.83 (+108.88%)	1.69 (-11.	.41%) 3.48 (+101.84%)	1.45 (-4.87%)	4.13 (+116.57%)	
BerkeleyLM	H.3	(+32.90%) 6.70 (+258.81%)	(+72.70%) <b>0.97</b> (-28.46%)	(+22. <b>7.82</b> (+310.	.04%) (+61.70%) .38%) <b>1.13</b> (-34.35%)	(+10.83%) 9.24 (+507.79%)	(+79.76%) <b>2.18</b> (+13.95%)	
BerkeleyLM	H.50	(+423.40%) <b>7.96</b> (+326.03%)	(-40.85%) <b>0.97</b> (-28.49%)	(+465. 9.37 (+391.	.36%) (-47.41%) .32%) <b>0.96</b> (-44.27%)	(+608.07%)	(-5.42%)	
Expgram		$2.06^{(+521.45\%)}_{(+10.18\%)}$	(-40.88%) <b>2.80</b> (+106.61%)	<b>2.24</b> (+17.	.87%) (-55.35%) .36%) <b>9.23</b> (+435.33%)	_	_	
KenLM T.		2.99 (263X <sup>6</sup> ) (+133.56%)	(-5.47%)	3.44 285	52%) <b>1.94</b> (+12.32%)	—		
Marisa		<b>3.61</b> (+93.09%) (+181.66%)	<b>2.06</b> (+52.00%) (+25.67%)	<b>3.81</b> (+99. (+174.	(-10.01%) (-10.01%) (.60%) <b>3.24</b> (+87.96%) (+50.58%) (+50.58%)	—	_	
RandLM		1.81 (-3.06%)	4.39 (+224.20%)	2.02 (+6.	18%) <b>5.08</b> (+194.35%)	2.60 (+70.73%)	9.25 (+384.54%)	
		(+41.41%)	(+168.04%)	(+46.	.29%) (+135.82%)	(+98.90%)	(+302.19%)	
		Europarl		YahooV2			GoogleV2	
--------------	-----------	---------------------------------------	------------------------------------	-------------------------------	---	-------------------	-------------------------------------	--
	bpg	1	$us \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$	
PEF-Trie	1.87	[]	1.35	1.91	1.73	1.52	1.91	
PEF-RTrie	1.28		1.64	1.38	2.15	1.31	2.30	
BerkeleyLM (	C. 1.70	(-8.89%)	2.83 (+2°X <sup>8%)</sup>	1.69 (-11.4	41%) <b>3.48</b> (+ <b>2</b> × <sup>4%)</sup>	<b>1.45</b> (-4.	87%) <b>4.13</b> (+1 <b>2X</b> %)	
BerkeleyLM H	1.3 6.70	(+258,81%)	(+72.70%) 0.97 (-28.46%)	$7.82_{(+310.3)}$	(+61.70%) $(+61.70%)$ $(-34.35%)$	9.24 (+507.	$\frac{83\%}{20\%} 2.18 (+13.95\%)$	
BerkeleyLM H	1.50 7.96	(-42ao310%) (-4522X <sup>6</sup> )	(-40.85%) <b>0.97</b> (-28.49%)	9.37 (- <b>5</b> -8	(-47.41%) (-47.41%) (-44.27%)	(+608.	(-5.42%)	
Expgram	2.06	(+521.45%) (+10.18%)	(-40.88%) 2.80 (+126 (21%)	(+576.) <b>2.24</b> (+17.)	$\begin{array}{c} (-55.35\%) \\ (-35.35\%) \\ 9.23 \\ (-3355\%) \\ 3.55 \\ \end{array}$	<b>)</b> –		
KenLM T.	2.99	(2:3X°)	1.28 (-5.47%)	<sup>3.44</sup> 2.5	(1932-270) (1.94) (+12.32%)	5 –		
Marisa	3.61	(+133.56%) (+93.09%)	(-21.84%) <b>2.06</b> (+52.00%)	3.81 (+148.)	52%) (-10.01%) 50%) <b>3.24</b> (+87.96%)	_		
		(+181.00%)	(+25.67%)	(+174.9	93%) (+50.58%)	)		
RandLM	1.81	(-3.06%)	4.39 (-2252%)	2.02 (+6.)	18%) <b>5.08</b> (2194.37%)	<b>2.60</b> (+70.	73%) 9.25 (+36.54%)	
		(+41.41%)	(+168.04%)	(+46.2	29%) (+135.82%)	(+98.	90%) (+302.19%)	

		Europarl		YahooV2	GoogleV2	
	bpg	$\mu s \times query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie PEF-RTrie	$1.87 \\ 1.28$	1.35 $1.64$	$1.91 \\ 1.38$	$     \begin{array}{c}       1.73 \\       2.15     \end{array} $	$1.52 \\ 1.31$	$1.91 \\ 2.30$
BerkeleyLM C.	1.70 (-8.	$\begin{array}{c} 2.83 \ (+12878\%) \\ (+72.70\%) \\ (+72.70\%) \end{array}$	<b>1.69</b> (-11	(41%) 3.48 (+2×4%)	1.45 (-4.	$\begin{array}{c} 4.13 (+12 \times \%) \\ (+79.76\%) \end{array}$
BerkeleyLM H.3	6.70 (+258,	(+12.10%) 81%) <b>0.97</b> (-28.46%)	7.82 (+310	(-34.35%) <b>1.13</b> (-34.35%)	9.24 (+507	<b>2.18</b> (+13.95%)
BerkeleyLM H.50	7.96 (45-2	(-40.85%) = (-40	9.37 (-5-1	(-44.27%) (-44.27%) (-55.25%)		
Expgram	<b>2.06</b> (+10.	$\begin{array}{c} (-40.38\%) \\ (18\%) \\ (2.80 \ (+186.61\%) \\ (170.80\%) \\ (170.80\%) \end{array}$	2.24 (+17	.36%) 9.23 (-3355%)	—	
KenLM T.	2.99 2.0	(-5.47%)	3.44 2.4	<b>5X</b> <sup>()</sup> <b>1.94</b> (+12.32%)	) –	
Marisa	3.61 (+133)	$\begin{array}{c} (-21.84\%) \\ (09\%) \\ (09\%) \\ (+52.00\%) \\ (+25.67\%) \end{array}$	(3.81 (+99))	$ \begin{array}{c} (-10.01\%) \\ (60\%) \\ (.98\%) \\ \end{array} \begin{array}{c} 3.24 \\ (+87.96\%) \\ (+50.58\%) \end{array} $	—	_
RandLM	1.81 (-3. (+41.	$\binom{0.6\%}{4.1\%} \begin{pmatrix} 4.39 & (2^{12}53\%) \\ (+168.04\%) \end{pmatrix}$	<b>2.02</b> (+6 (+46	.18%) <b>5.08</b> (2 <sup>19</sup> <b>5.0%</b> ) (+135.82%)	<b>2.60</b> (+70. (+98.	73%) <b>9.25</b> (+335,54%) 90%) (+302.19%)

		Europarl		YahooV2	GoogleV2	
	bpg	$\mu s  imes query$	bpg	$\mu s \times query$	bpg	$\mu s \times query$
PEF-Trie PEF-RTrie	$1.87 \\ 1.28$	$\begin{array}{c}1.35\\1.64\end{array}$	$1.91 \\ 1.38$	$     \begin{array}{c}       1.73 \\       2.15     \end{array} $	$1.52 \\ 1.31$	$1.91 \\ 2.30$
BerkeleyLM C.	1.70 (-8	(+72.70%)	<b>1.69</b> (-11 (+22	(41%) 3.48 (+2×4%)	1.45 (-4.	$\begin{array}{c} 4.13 (+12 \mathbf{X}^{\%}) \\ \mathbf{4.3\%} & (+79.76\%) \end{array}$
BerkeleyLM H.3	$6.70_{(+258)}$	(-28.46%) (0.97 (-28.46%)	$7.82_{(+310)}$	(-34.35%) <b>1.13</b> (-34.35%) (-47.41\%)	$9.24_{(+507)}_{(+608)}$	2.18 (+13.95%) (-5.42%) (-5.42%)
BerkeleyLM H.50	7.96 (+52)	<b>0.97</b> (-28.49%) (-40.88%)	9.37 (-59)	<b>0.96</b> (-44.27%) .87%) (-55.35%)	_	_
Expgram	2.06 (+10	$\begin{array}{c} (13\%) \\$	2.24 (+17	.36%) 9.23 (-3355%)	—	_
KenLM T.	2.99 (260) (+133)	(-21.84%)	<b>3.44</b> (2 <sup>80</sup> )	<b>1.94</b> (+12.32%) (-10.01%)	) –	
Marisa	3.61 (+93) (+73) (+73)	$\left(\begin{array}{c} 0.0\%\\ 0.0\%\end{array}\right) \begin{array}{c} 2.06 & (+52.00\%)\\ (+25.67\%) & (+25.67\%) \end{array}$	$(3.81)^{(+99)}_{(+114)}$	60%) <b>3.24</b> (+87.96%) (+50.58%) (+50.58%)	—	—
RandLM	1.81 (-3 (+41	.06%) <b>4.39</b> (************************************	<b>2.02</b> (+6 (+46	.18%) <b>5.08</b> (2 <sup>19</sup> <b>5.0%</b> ) (+135.82%)	<b>2.60</b> (+70. (+98.	73%) <b>9.25</b> (+37) (43) 90%) (+302.19%)

- Elias-Fano Tries substantially **outperform ALL** previous solutions in **both space and time**.
- As fast as the state-of-the-art (KenLM) but more than twice smaller.

#### Scalable Modified Kneser-Ney Language Model Estimation





seconds

1.3 GB 233,035,325 total words

> 1,255,027 20,431,391 82,815,629 153,984,231 196,779,246

> 455,265,524

3.2 GB 495,527,349 total words

> 15,039,323 44,033,774 142,894,817 280,714,113 381,284,741 \_\_\_\_\_\_ 863,966,768

#### Scalable Modified Kneser-Ney Language Model Estimation

seconds







3.2 GB 495,527,349 total words

> 15,039,323 44,033,774 142,894,817 280,714,113 381,284,741 863,966,768



#### Scalable Modified Kneser-Ney Language Model Estimation

seconds







3.2 GB 495,527,349 total words

15,039,323
44,033,774
142,894,817
280,714,113
381,284,741
863,966,768



seconds

seconds

#### Scalable Modified Kneser-Ney Language Model Estimation

seconds

Tongrams - 1

32 22











863,966,768



counting normalization interpolation

seconds

#### Scalable Modified Kneser-Ney Language Model Estimation













seconds



counting normalization interpolation

seconds





### 2. Develop other research ideas



✓ Inverted Indexes with *false positives* allowed.

2. Develop other research ideas







3. 6 months abroad.

# Thanks for your attention, time, patience!

Any questions?