

CSC_51063 Final : Prediction and Entropy of Printed English (German, French, etc)

Prediction and Entropy of Printed English

By C. E. SHANNON

(Manuscript Received Sept. 15, 1950)

A new method of estimating the entropy and redundancy of a language is described. This method exploits the knowledge of the language statistics possessed by those who speak the language, and depends on experimental results in prediction of the next letter when the preceding text is known. Results of experiments in prediction are given, and some properties of an ideal predictor are developed.

1. INTRODUCTION

IN A previous paper¹ the entropy and redundancy of a language have been defined. The entropy is a statistical parameter which measures, in a certain sense, how much information is produced on the average for each letter of a text in the language. If the language is translated into binary digits (0 or 1) in the most efficient way, the entropy H is the average number of binary digits required per letter of the original language. The redundancy, on the other hand, measures the amount of constraint imposed on a text in the language due to its statistical structure, e.g., in English the high frequency of the letter E , the strong tendency of H to follow T or of U to follow Q . It was estimated that when statistical effects extending over not more than eight letters are considered the entropy is roughly 2.3 bits per letter, the redundancy about 50 per cent.

Since then a new method has been found for estimating these quantities, which is more sensitive and takes account of long range statistics, influences extending over phrases, sentences, etc. This method is based on a study of the predictability of English; how well can the next letter of a text be predicted when the preceding N letters are known. The results of some experiments in prediction will be given, and a theoretical analysis of some of the properties of ideal prediction. By combining the experimental and theoretical results it is possible to estimate upper and lower bounds for the entropy and redundancy. From this analysis it appears that, in ordinary literary English, the long range statistical effects (up to 100 letters) reduce the entropy to something of the order of one bit per letter, with a corresponding redundancy of roughly 75%. The redundancy may be still higher when structure extending over paragraphs, chapters, etc. is included. However, as the lengths involved are increased, the parameters in question become more

¹ C. E. Shannon, "A Mathematical Theory of Communication," *Bell System Technical Journal*, v. 27, pp. 379-423, 623-656, July, October, 1948.

Table of Contents

1. Motivation
2. Entropy, Entropy Rate and Redundancy
3. Character-level Entropy
4. Word-level Entropy and Zipf's Law
5. Prediction, Illustration, and Limits
6. Conclusion

Problem Framing

Lecture 3 Example

THE ENGLISH LANGUAGE(I)

- Approximation of order 0 (all symbols, don't forget the "space", are i.i.d.):

XFOML RXKHRJFFJUJ ZLPWCFWKCYJ FFJEYVKCQSGXYD QPAAMKBZAACIBZLHJQD

$$H_0 = \log_2 27 \approx 4.76$$

- Approximation of order 1 (the letters are chosen according to their frequency in English)

OCRO HLI RGWR NMIELWis EU LL NBNESEBYA TH EEI ALHENHTTPA OOBTTVA NAH BRL

$$H_1 \approx 4.03$$

- Approximation of order 2 : same distribution of the pairs as in English

ON IE ANTSOUTINYS ARE T INCTORE ST BE S DEAMMY ACHIN D ILONASIVE TUCOOWE AT
TEASONARE FUSO TIZIN ANDY TOBE SEACE CTISBE

- Approximation of order 3 : same frequency of the triplets as in English

IN NO IST LAT WHEY CRATICT FROURE BERS GROCID PONDENOME OF DEMONSTURES OF THE
REPTAGIN IS REGOACTIONA OF CRE

Any idea to generate correctly some English test?

2) Entropy, Entropy Rate and Redundancy

N-gram entropies from text

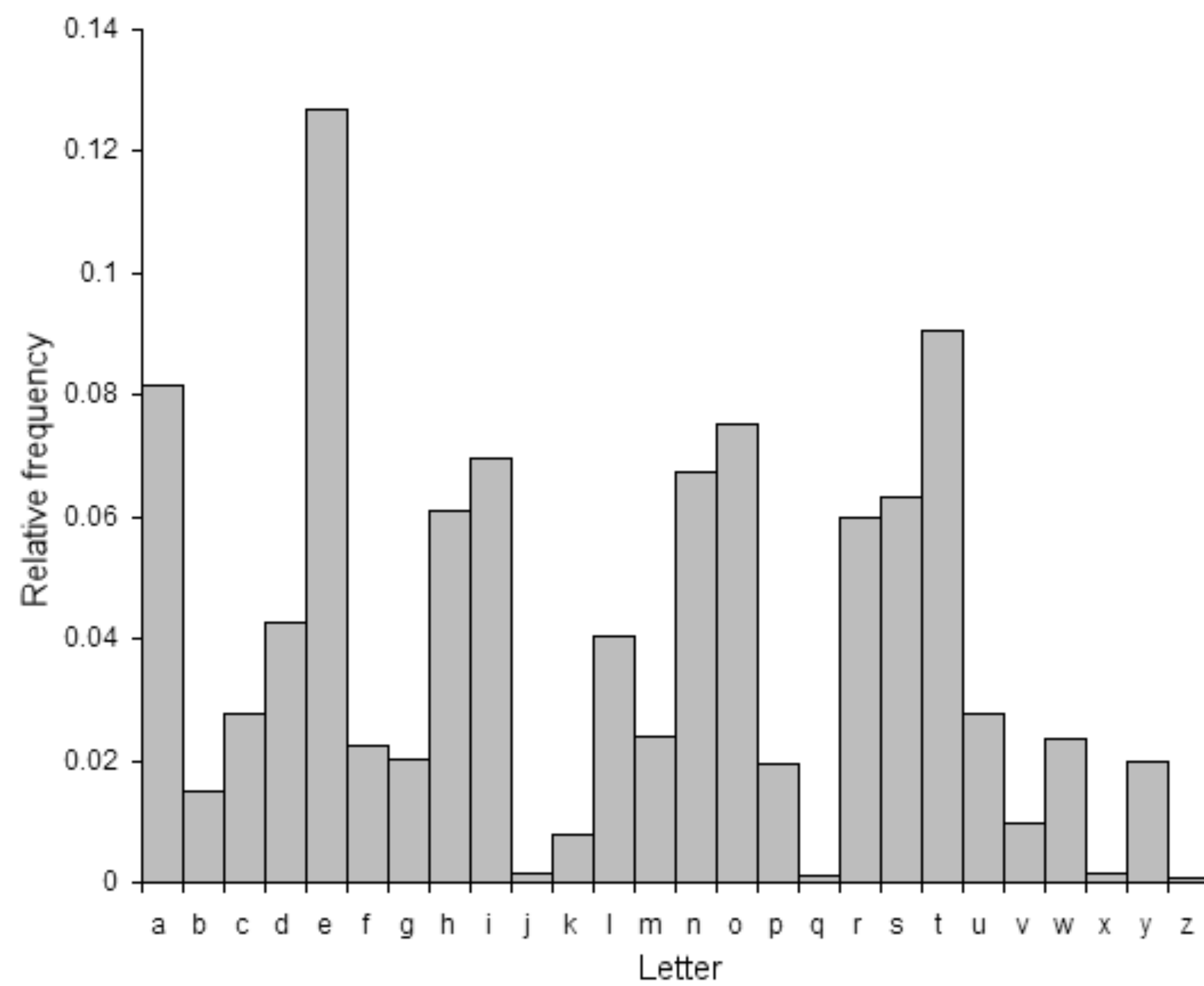
Introduction and Definitions

- Zero Order Model
- Maximum entropy (uniform letters) = 4.75 bits per letter
- given 26 letters + whitespace and assuming that all characters are equally likely

- First Order Model takes character frequencies into account

$$H_1 = H(X_t) = - \sum_x p(x) \log p(x)$$

- However this is assuming zero memory, ie characters appear independently from what was written before



Letter and next-letter frequencies in English
measured across 1 million articles from Wikipedia

	t ^h	a ⁿ	o ^r	s ^t	i ⁿ	c ^o	w ^a	b ^e	p ^r	f ^o	m ^a	h ^a	r ^e	d ^e		e ^x	l ^a	n ^o	g ^r	u ⁿ	v ⁱ	j ^o	k ⁱ	y ^e	q ^u	z ^e	x
e		r	n ^t	d	s	a ^r	l ^e	c ^t	t	m ^e	e ⁿ	v ^e	i ^r	p ^t	x ^p	g ⁱ	f ^e	w	y	o ^p	b ^r	h ^o	u ^r	q ^u	k	z	j ^e
a	n ^d	t ⁱ	r ^e	l	s		c ^t	m ^e	d	i ⁿ	g ^e	y	p ^p	b ^l	v ^e	u ^s	k ^e	f ^t	w ^a	h	e ^l	z ⁱ	j ^o	x	a	o	q ^u
t	h ^e		e ^r	i ^o	o	a ^t	r ^a	s	u ^r	y	t ^e	l ^e	w ^o	c ^h	m ^e	b ^a	n ^e	f ^o	v	z	p ^u	d ^o	g ^r	k ⁱ	j ^a	x	q ^u
i	n	s	t ^h	o ⁿ	c ^a	l ⁱ	e ^s	a ⁿ	r	v ^e	m ^e	d ^e	g ^h	f ⁱ		p	b ^e	z ^e	k ^e	u ^m	x	i	q ^u	j ⁱ	h ^a	w ^a	y ^a
n		d	t	g	e	s	a ^l	i ^t	c ^e	o ^t	n ^e	u ^m	y	f ^o	l ^y	k	v ^e	m ^e	b ^s	h ^a	r ^e	w ^a	j ^u	p ^r	z ^e	q ^u	x
o	n	r	f		u ⁿ	m	l ⁱ	t ^h	w ⁿ	s ^t	p ^e	v ^e	o ^k	d ^e	c ^a	g ^r	b ^e	i ⁿ	a ^d	k	y ^a	e ^s	h ⁿ	x ⁱ	j ^e	z ^e	q ^u
s		t	e	i ⁿ	o ⁿ	s	h ^e	u ^c	a ⁿ	p ^e	c ^h	m ^a	l ^a	y ^s	k ⁱ	w ⁱ	b ^a	f ^u	n ^e	r ^a	d ^a	q ^u	g ^o	v ⁱ	j ^o	z ^e	x
r	e		i ⁿ	o ^m	a ^t	s	t	y	d	n	m ^a	c ^h	u ^c	r ⁱ	g ^e	l ^y	k	v ⁱ	p ^o	b ^a	f ^o	h ^o	w ^a	z ^e	q ^u	j ^a	x ⁱ
h	e		a ^t	i ^s	o ^u	t	r ^o	u ^r	y	n	s	l ^y	m ^e	w ^a	b ^o	d ^r	c ^o	f ^u	p	h ^e	k ^e	v ⁱ	g ^a	q ^u	z	j ^o	x
l	e		i ⁿ	a ⁿ	l ^o	w ^y	d	s	u ^d	t ^h	m	f	v ^e	b ^u	p ^h	k	c ^o	w ^a	g ^a	r ^e	h ^e	n ^e	j ^a	z ^a	q ^a	x	
d		e ^r	i ⁿ	a ^y	o ⁿ	u ^c	s	r ^e	d ⁱ	y	l ^e	g ^e	t ^h	m ⁱ	v ^a	w ^a	h ^e	c ^a	n ^e	b ^e	f ^o	p	j ^a	q ^u	k ⁱ	z ⁱ	x
c	o ⁿ	h	e	a ^l	t ⁱ	i ^a		l ^u	r ^e	k	u ^l	c ^e	s	y	d	q ^u	m	b ^s	p	n ^e	z ^e	f	g ^r	w	v	j ^a	x
u	s ^e	n ^d	r ^e	t	l ^a	c ^h	m	e	a ^l	g ^h	p	d ^e	i ^l	b ^l		k	f ^f	o ^u	x	v ^e	y	z ^z	h ^a	j ⁱ	u ^m	w ^a	q ^u
m	e	a ⁿ		o ⁿ	i ⁿ	p ^l	b ^e	u ^s	m ^e	s	y	t ^h	c ^c	n	d ^a	l ^e	r	h ^e	f ^o	w ^e	g	k	v	j ^a	x	z ^a	q ^u
f		o ^r	i ^r	r ^o	e ^r	a ^c	f ^e	t ^e	u ^l	l ^o	s	y	c	m ^a	g ^h	d ^e	b ⁱ	p ^r	h ^e	w ^a	n ^e	j ^o	k	x	v ⁱ	z ^e	q ^u
p	e ^r	r ^o	o ^r	a ^r	l ^a		i ⁿ	u ^b	p ^e	h ⁱ	t ⁱ	s	m ^e	y	k ^m	c	d ^a	n ^b	g ^r	b ^e	f	w ^a	v	j ^a	z ⁱ	x	q
g		e	h ^t	a ⁿ	r ^a	i ⁿ	o ^v	u ^a	n ^e	l ^e	s	y	t ^h	g ^e	d ^o	m ^e	b ^y	w ^r	f ^o	p ^o	c ^o	k ^o	j ⁱ	z ^h	v ⁱ	x ⁱ	q ⁱ
w	a ^s	i ^t	e ^r	h ⁱ	o ^r		n	s	r ⁱ	l ^e	t ^h	y ^e	d ^e	w ^e	b ^o	k	c ^a	f ^o	m ^a	u	p ^o	g ^r	v	j ^a	z ^a	x	q ^u
y		e ^a	s	o ^u	a ^l	t ^h	i ⁿ	p ^e	l ^e	m ^p	n ^a	c ^l	r ⁱ	d ^r	b ^o	w ^h	h ^e	u ^g	f ^o	g ^e	z ^a	k ^e	v ⁱ	j ^o	y ^a	x	q ^u
b	e ^r	a ⁿ	o ^u	y	l ^e	u ^t	r ⁱ	i ^l	s ^p		b ^e	j ^e	c	t ^a	m ^a	n	d ⁱ	h ^a	v ⁱ	f	w ^a	p ^a	g ^r	k	z ^e	x	q
v	e ^r	i ⁿ	a ^l	o ^l		y	s	u ^l	d	r ^e	l ^a	t ^h	c	n	p	g	h ^s	m	f	b	v ⁱ	k ^a	w	x	j	z	q
k		e	i ⁿ	s	a	n ^o	o ^r	l ^y	y	h ^a	u ^r	r ^a	t ^h	m	w ^a	b ^u	g ^r	k ^a	f ^o	p ^o	d ^o	c ^o	v ⁱ	j ^a	x	z	q ^u
x		p ^e	t ^e	i ^s	a ^m	e ^c	c ^e	u ^a	o ⁿ	h ⁱ	y	f ^o	m ^e	s	x	l ^e	b ^o	v ⁱ	r ^a	w ^e	n ^e	d ^e	g ^a	j	k	q ^u	z
j	o ^r	u ^s	a ⁿ	e ^c	i ^m		r	p ^g	s	k	d	n ^a	t ^h	c	m	h ^a	j	l ⁱ	f ^k	b	v ^o	y ^a	g	w	z	x	q
z	e ^d	a ^t	i ⁿ		o ⁿ	z	u ^r	h ^a	y	l ^e	s	t ^h	b ^e	m ^a	n ^a	k ^o	d ^a	w ^e	r ^a	g ^e	c ^z	v ^o	f ^e	p ^a	q ^u	x	j ^a
q	u ^e		i	a ^e	s	t ^h	r ^t	l	o ^m	b ^a	e	n ^b	c	q ^a	v ⁱ	f	h	m	p	w ^e	d ^a	x	g ^e	z	j ⁱ	k	y ^t

Colour key

a Letters and spaces sorted by overall frequency. Ignores case and accents.

n^d Next letter sorted by frequency. Small letter is the most common third letter following the pair.

Blank letters indicate spaces

Letter frequencies

		31.6% to 100.0%
		10.0% to 31.6%
		3.2% to 10.0%
		1.0% to 3.2%
		0.3% to 1.0%
		0.1% to 0.3%
		0.0% to 0.1%

Markov generators

The letters distributions in the chart can be used to generate psuedowords such as the ones below. A similar approach, at the word level, is used for online parodv generators.

bastrabot	dithely
loctrion	raliket
calpereek	amorth
forliatitive	asocult
wasions	quarm
felogy	winterlitterand
sonsih	uniso
fourn	hise
meembege	nuouish
prouning	guncelawits
nown	rectere
abrip	doesium

Letter and next-letter frequencies in French

measured across a selection of articles from Wikipedia

	d ^e	l ^e	e ^t	p ^a		a	c ^o	s ^e	m ^a	t ^r	r ^e	f ^o	u ⁿ	b ^a	i ⁿ	n ^o	v ⁱ	o ^u	g ^r	z ^o	q ^u	h ^a	j ^u	k ^a	w ⁱ	y	x ⁱ
e		s	n ^t	r	t	l ⁱ	m ^e	u ^r	c ^t	e	p ^o	a ^u	d ⁱ	g ⁱ	v ^e	i ⁿ	z	x ^t	f ^e	b ^a	o ^r	q ^u	y	h ⁱ	j ^o	k	w
a		n ^t	r ^t	i ^s	l ^e	t ⁱ	u	s ^s	g ^e	c ^e	m ^e	v ^e	p ^p	b ^l	d ^e	y ^s	e	f ^f	o	z ^z	k ^a	h ^a	a	q ^u	x ⁱ	j ^o	w ^a
i	n ^e	e	s	t ^e	o ⁿ	l		c ^a	r ^e	a ⁿ	q ^u	m ^e	g ⁿ	d ^e	v ^e	p ^e	f ⁱ	b ^l	k ⁱ	i	x	z ^z	u ^m	j ^a	y ^a	h ^a	w ^a
n		t	e	s	a ⁱ	c ^e	o ^m	d ^e	i ^s	n ^e	g	u ^e	f ^a	v ⁱ	y ^m	z ^a	r ^e	l ^e	k	p ^a	q ^u	h ^a	j ^a	m ^a	b ^a	w ^a	x ⁱ
s		e	t	i ^t	s ^e	o ⁿ	a ⁿ	u ^r	p ^e	c ^o	h ^a	l ^a	m ^e	q ^u	y ^s	k ⁱ	b ⁿ	d ^e	z ^a	f ^a	r ^e	n ^o	v ^o	g ^e	w ⁱ	j ^e	x ⁱ
r	e		a ⁿ	i ^e	o ^u	t ⁱ	s	n ^a	d	r ^e	m ^e	c ^e	u ^s	g ^e	b ^e	v ^e	l ^e	p ^r	f ⁱ	y	k	q ^u	z ^e	h ^o	j ^a	w ⁱ	x ⁱ
t		e	i ^o	a ⁱ	r ^e	o ^u	s	u ^r	t ^e	h ^e	c ^h	y ^p	l ^e	z	p ^w	b ^a	m ^a	v	d ^e	n ^a	g ^e	f ^r	w ⁱ	k ^o	x ^a	j ^e	q ^u
o	n	u ^r	r ^t	m ^m		i ^s	l ^o	s	p ^e	t ^e	c ^a	g ^r	d ^e	v ⁱ	b ^l	f ^f	o ^t	y ^e	a	e	w	k ^o	h ⁿ	z ^z	q ^u	x ⁱ	j ^e
l	e	a	i ^e	l ^e		o ⁿ	u ^s	s	t ⁱ	m ^a	d	h ⁱ	b ^u	y	p ^h	g ^a	c ^o	v ^a	q ^u	f ^e	z ^a	k ^o	n ^e	j ^e	r ^e	w ^a	x
u	r	e		n ^e	s	i	t	l ^a	x	v ^e	c ^h	a ⁿ	m ^e	p ^e	d ^e	b ^l	g ^e	j ^o	f ^f	o ⁿ	z ^e	k	u ⁿ	y ^a	h ^a	q ^u	w ^e
d	e	a ⁿ	u	i ^s	o ⁿ		r ^e	s	h ^o	m ⁱ	y ⁿ	d ^a	l ^e	j ^a	c ^o	v ^e	g ^e	n ⁱ	p ^e	z ^z	b ^a	f	w ^a	t	k ^r	x ⁱ	q ^u
c	o ⁿ	e	h ^e	a ^r	i ^e	t ⁱ	l ^e		u ^l	r ^e	c ^e	k	s	y ^c	q ^u	d	n	m	f	z ^e	p ^o	b ⁱ	v	g ^e	j ^a	x ^d	w
p	a ^r	o ^u	e ^r	r ^e	i ⁿ	l ^u	h ⁱ	u ⁱ	p ^e	t ^e		s	f ^a	y ^r	c ^c	d ^f	w ^w	z ^z	g ^a	m ^s	n ^e	b ^e	x ^p	v	k ⁱ	j ^e	q ^r
m	e ⁿ	a ⁿ	i ⁿ	o ⁿ	m ^e	p ^o		b ^r	u ⁿ	s	y ^s	c ^o	t ^e	l ^a	n ^e	g ^e	r	h	d ^e	z ^a	f ^o	v ⁱ	w ^a	k ⁱ	j	q ^u	x
g	e	r ^a	a ^l	i ^o	n ^e	u ^e	o ^u		l ^e	h ^t	g ⁱ	s	y ^p	m ^e	t ^e	p ^o	d ^a	z ^h	b ^o	c ^e	f ^r	j ⁱ	w ^a	x ⁱ	v ⁱ	k ^o	q ⁱ
v	e ^r	i ^l	a ⁿ	o ⁱ	r ^e		u ^e	s ^k	l ^a	d	y	n ⁱ	c ^e	k ^o	p	g ^o	f	b ^e	h ^s	z ^v	t ^o	m ⁱ	v ^l	w ^x	j ^e	x ^u	q ^u
h	e	a ⁿ	i ^e	o ^m		u ⁱ	y ^s	r ⁱ	n ⁱ	t ^t	l ^e	s	m ^e	c	z ^z	p	w ^a	b ^a	d ^a	k ^o	v ⁱ	g ^e	h ^a	f ⁱ	j ^a	q ⁱ	x ^a
b	a ^s	r ^e	e ^r	i ^e	l ^e	o ^u	u ^t		s ^e	n	y	d ⁱ	b ⁱ	k ^m	j ^e	t ^e	c	w ^e	h ^a	p ⁱ	m ^a	z ^a	v ^e	f ^o	g ^e	x	q
f	i ^c	o ^r	e ^r	a ⁱ	r ^a	f ^e		u ^t	l ^e	s	t	c	g ^h	m ^a	p ⁱ	k	v ⁱ	d ^e	n ^a	y	z	b ⁱ	h ^e	x	j ^a	w ⁱ	q ^u
z	a ⁿ	o ⁿ		e ⁿ	i ⁿ	h ^a	z ^a	u ^r	y ^g	l ^a	w ⁱ	d ^e	v ^e	s ⁱ	b ⁱ	t ^e	m ^a	r ⁱ	n ^a	p ⁱ	k ^o	c ^e	g ^e	q ^a	f ^s	x	j ^a
q	u ^e		i ⁿ	a ^a	l ^e	r ^s	e ^y	x ^x	o ^u	q ^a	c ^o	h	j	s	v ⁱ	n ^e	b	f	w ^e	t ^a	d	z ^o	m	p ^o	y ^a	g ^z	k
y		s	a ⁿ	e ⁿ	m ^e	p ^e	o ^u	l ^e	n ^a	r ⁱ	c ^l	t ^h	d ^r	u	i ⁿ	g ^a	z ^o	b ^r	k ^l	v ^e	x	w ^o	h ^a	f ⁱ	j ^a	y ^a	q ^u
x		i ^s	e	t ^e	p ^l	a ⁿ	o ⁿ	c ^e	v ⁱ	u ^e	x ⁱ	y ^l	l ^e	z ^z	m ^e	s ⁱ	h ^u	f ^o	d ^e	g ^o	n ^o	q ^u	b ^o	r ^o	k ^e	j ^a	w ⁱ
k	a		i ⁿ	o ^v	e ^r	m	u ⁿ	h ^a	r ^a	s	l ^a	y	g	t ⁱ	k	n ^o	w ^a	p ^o	v	c ^e	z ^e	b ^o	d ^e	j ⁱ	f ⁱ	x	q ^u
j	o ^u	e ^u	a ⁿ	u ^s	i ^a		c	n ⁱ	k ^l	r	s ^k	p ^e	m ^o	d ^a	l ^e	v ^a	b ^a	j ⁱ	t ^a	h ^o	f ^a	y	z ^e	g ^o	w	x ^t	q ^t
w	i ^k	a ^r	e ^b		o ^r	w ^w	h ⁱ	s	u	n	y ^a	z ^e	l ^a	r ⁱ	p ⁱ	c	t ^a	m ^a	k ^a	d ^e	b ^o	q	x ^y	f ^o	g ^c	v ⁱ	j ^o

Colour key

a	Letters and spaces sorted by overall frequency. Ignores case and accents.
n ^d	Next letter sorted by frequency. Small letter is the most common third letter following the pair.

Blank letters indicate spaces

Letter frequencies

		31.6% to 100.0%
		10.0% to 31.6%
		3.2% to 10.0%
		1.0% to 3.2%
		0.3% to 1.0%
		0.1% to 0.3%
		0.0% to 0.1%

Markov generators

The letters distributions in the chart can be used to generate psuedowords such as the ones below. A similar approach, at the word level, is used for online parodv generators.

<i>cillesil</i>	<i>sulskini</i>
<i>lidhemin</i>	<i>plumeme</i>
<i>bachogine</i>	<i>crout</i>
<i>taphie</i>	<i>provicas</i>
<i>copit</i>	<i>odzzaccet</i>
<i>extreiles</i>	<i>pipiphien</i>
<i>chetratagne</i>	<i>outif</i>
<i>suro</i>	<i>extellages</i>
<i>nans</i>	<i>nutopune</i>
<i>entote</i>	<i>sporese</i>
<i>zhiquis</i>	<i>edes</i>
<i>aliet</i>	<i>randamelle</i>

N-gram entropies from text

Introduction and Definitions

- we now take **context** into account and list the most common digrams in english
- these are in order : th, he, in, en, nt, re, er, an, ti, es, on, at, se, nd ...
- this is where conditional entropy comes into play !
- from the digrams we obtain $P(X_t | X_{t-1})$
- we can define $H_2 = H(X_t | X_{t-1}) = \sum_{x,y} p(x,y) [-\log p(y | x)]$
- Shannon calls this 2nd order entropy (1st order Markov chain)

N-gram entropies from text

Introduction and Definitions

- but this is still not enough, English is not 1st order Markov chain
- ‘tio’ almost always forces ‘n’ where the context required is length 3
- also word-level and semantic constraints extend far beyond characters
- Shannon’s idea = if we don’t know the true dependency length -> approximate it with longer and longer finite contexts
- and the “true” entropy rate will be

$$H = \lim_{N \rightarrow \infty} H_N$$

- where H_N = conditional entropy of the next symbol given the previous N-1

$$H_N = H\left(X_t \mid X_{t-1}, \dots, X_{t-(N-1)}\right) = \sum_{x_1, \dots, x_N} p(x_1, \dots, x_N) \left[-\log p(x_N \mid x_1, \dots, x_{N-1})\right]$$

N-gram entropies from text

Introduction and Definitions

- Redundancy is the gap between Entropy and Encoding Efficiency
- Redundancy $R = 1 - \frac{H}{H_{\max}}$ where $H_{\max} = \log_2 |\mathcal{A}|$
- Redundancy means you can delete or corrupt part of the message and still recover it.
- eg : ths sntnc s stll rdbl -> this sentence is still readable

3) Character-level Entropy

Experiment Setting

- 30 books from project Gutenberg downloaded as .txt files
- Clean and tokenized either as char or words
- Markov chain of orders $k=1,2,3 \dots$ fitted on the corpora
- Calculate entropy and generate text samples

English	Carroll - Alice's Adventures in Wonderland Fitzgerald - The Great Gatsby Hobbes - Leviathan Austen - Pride and Prejudice Shakespeare - Romeo and Juliet Machiavelli - The Prince Dickens - A Tale of Two Cities Tolstoy - War and Peace Dostoevsky - The Brothers Karamazov Dostoevsky - Crime and Punishment
French	Voltaire - Candide, ou l'Optimisme Leroux - Le Fantôme de l'Opéra Flaubert - Madame Bovary Dumas - Le Comte de Monte-Cristo Hugo - Notre-Dame de Paris Stendhal - Le Rouge et le Noir Dumas - Les Trois Mousquetaires Verne - Voyage au centre de la Terre Verne - L'Île mystérieuse Verne - De la Terre à la Lune Verne - Le Tour du monde en quatre-vingts jours Verne - Vingt mille lieues sous les mers
German	Kant - Der Wille zur Macht Kant - Kritik Hermann Hesse - Steppenwolf Struwwelpeter Kafka - Die Verwandlung Kafka - Das Urteil Mann - Zauberberg 1&2

N-gram entropies from text

Comparison with Result from Shannon

- Shannon also calculated N-gram entropies for small values of N
- Based on the book “Secret and Urgent” by Fletcher Pratt
- Unfortunately I couldn’t find a txt version online so i opted for another book by the same author
- The Blue Star by Fletcher Pratt written in 1950

The N -gram entropies F_N for small values of N can be calculated from standard tables of letter, digram and trigram frequencies.² If spaces and punctuation are ignored we have a twenty-six letter alphabet and F_0 may be taken (by definition) to be $\log_2 26$, or 4.7 bits per letter. F_1 involves letter frequencies and is given by

$$F_1 = - \sum_{i=1}^{26} p(i) \log_2 p(i) = 4.14 \text{ bits per letter.} \quad (3)$$

The digram approximation F_2 gives the result

$$\begin{aligned} F_2 &= - \sum_{i,j} p(i,j) \log_2 p(i,j) \\ &= - \sum_{i,j} p(i,j) \log_2 p(i,j) + \sum_i p(i) \log_2 p(i) \\ &= 7.70 - 4.14 = 3.56 \text{ bits per letter.} \end{aligned} \quad (4)$$

² Fletcher Pratt, “Secret and Urgent,” Blue Ribbon Books, 1942.

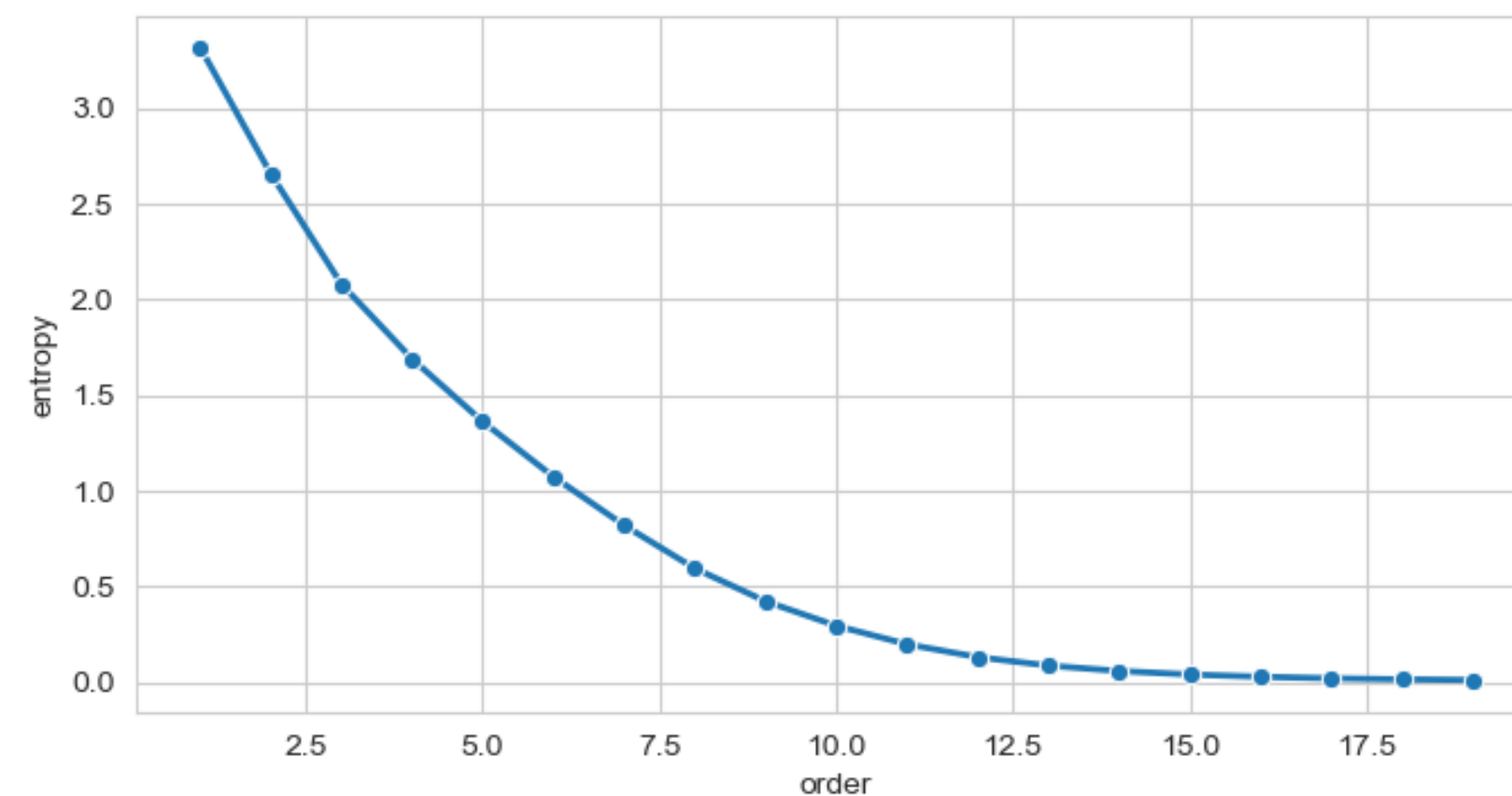
	F_0	F_1	F_2	F_3	F_{word}
26 letter.....	4.70	4.14	3.56	3.3	2.62
27 letter.....	4.76	4.03	3.32	3.1	2.14

Results from Shannon Paper

N-gram entropies from text

Empirical Estimates from Text

- Import and tokenize the text
- For all N-grams count distinct next step tokens and normalize to probability
- Calculate entropy



Alice's Adventures in Wonderland

CHAPTER I

Down the Rabbit-Hole

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had pee

Total tokens: 143407

Unique tokens: 31

['\n' ' ' ', ' . ' ? ' 'a' 'b' 'c' 'd' 'e' 'f' 'g' 'h' 'i' 'j' 'k' 'l' 'm' 'n' 'o' 'p' 'q' 'r' 's' 't' 'u' 'v' 'w' 'x' 'y' 'z']

----- Model Statistics -----

Order of N-gram Model (N) : 2
Computation Time (s) : 0.0241
Number of Unique Contexts: : 544
Total Observed N-grams (Transitions) : 143,405
Unique (Context, Token) N-grams : 4,211
Conditional Vocab Size (Next Tokens) : 31

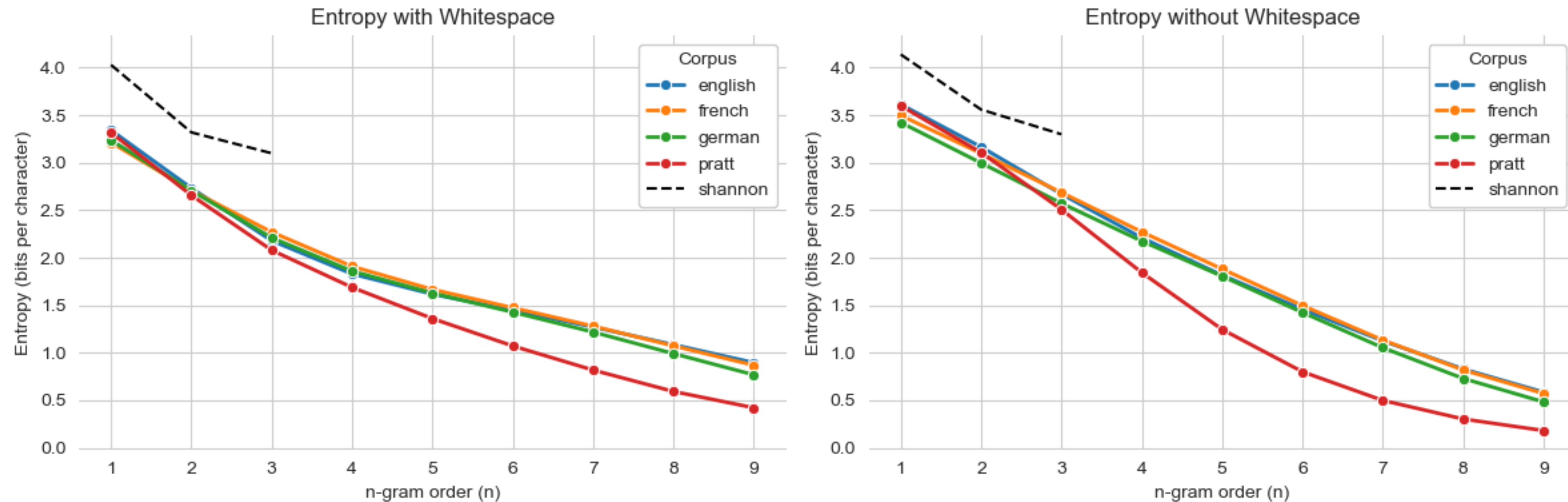
Most likely token after ('s', ' ') :

	token	proba
5	t	0.138748
0	a	0.131980
8	s	0.097011
16		0.077270
1	i	0.068246

Model Entropy: 2.5795 bits

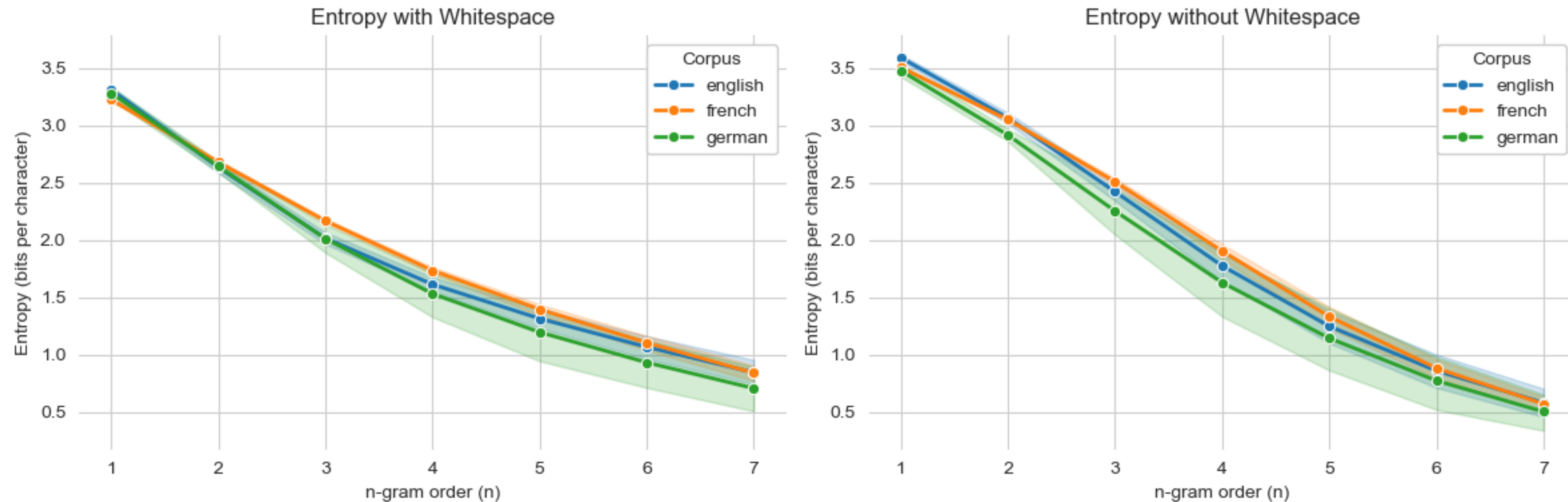
Char-level Entropy Across Languages

Comparison with tabulated values from Shannon



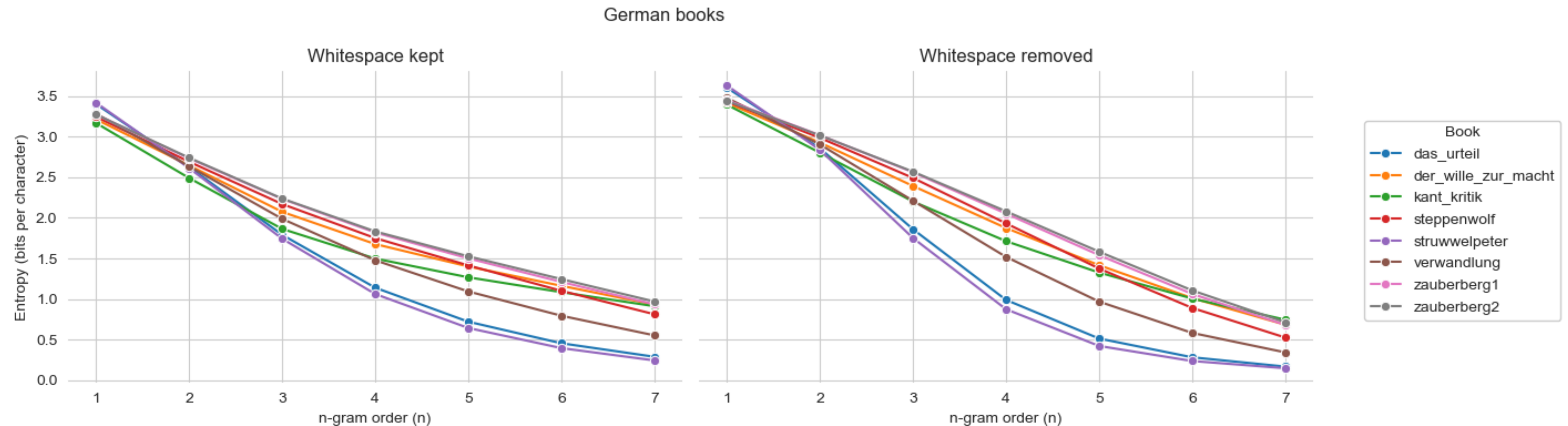
Char-level Entropy Across Languages

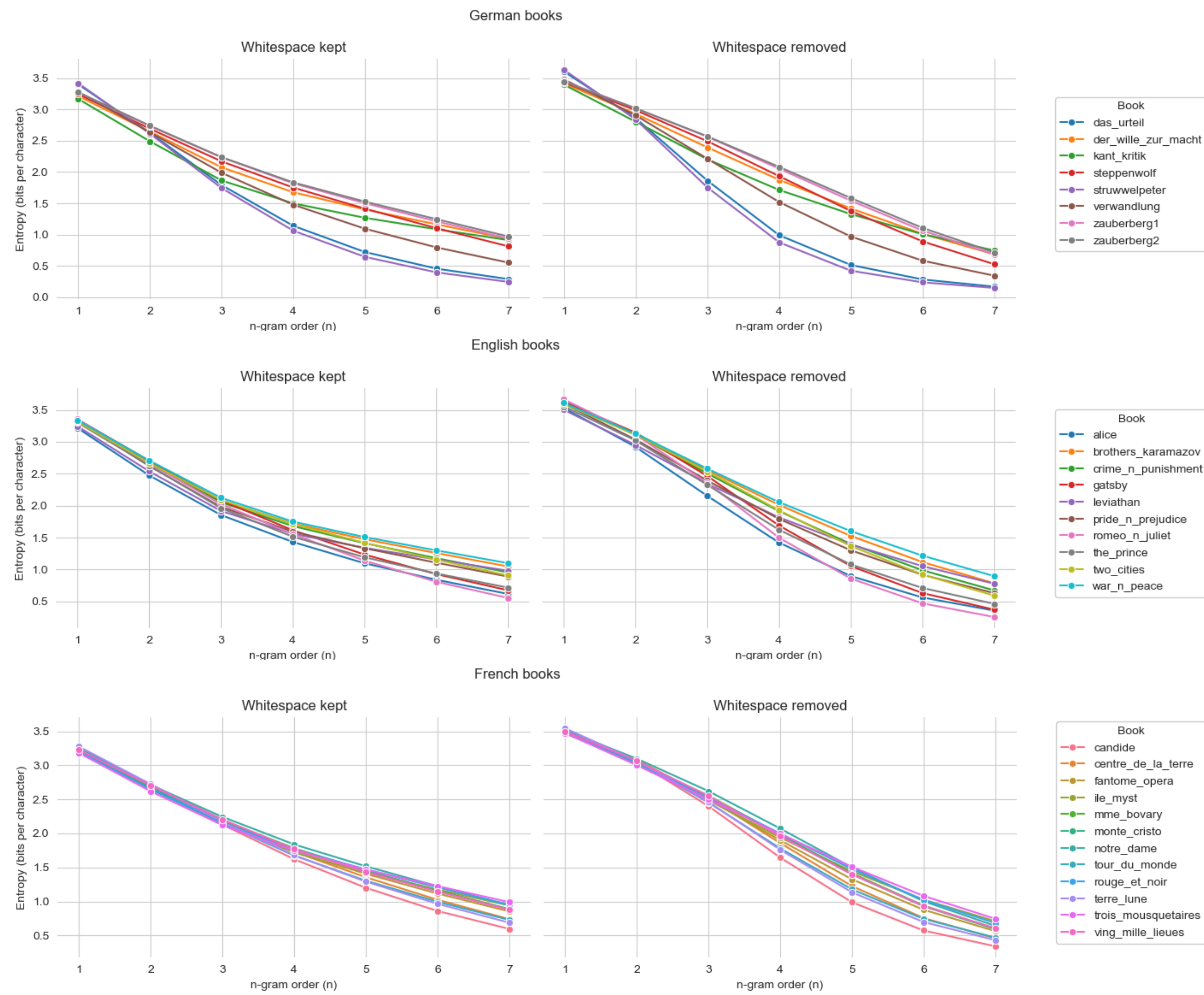
Book level analysis for each language



Char-level Entropy Across Languages

Book level analysis for each language





4) Word-level Entropy and Zipf's Law

N-gram entropies from text

Comparison with Result from Shannon

- Shannon also calculated word-level N-gram entropies for small values of N
- But only based on estimation from tabulated word frequencies and obtaining 11.82 bits per word
- and then by taking the average word length to be 4.5 letters
- $11.82 / 4.5 = 2.62$

The N -gram entropies F_N for small values of N can be calculated from standard tables of letter, digram and trigram frequencies.² If spaces and punctuation are ignored we have a twenty-six letter alphabet and F_0 may be taken (by definition) to be $\log_2 26$, or 4.7 bits per letter. F_1 involves letter frequencies and is given by

$$F_1 = - \sum_{i=1}^{26} p(i) \log_2 p(i) = 4.14 \text{ bits per letter.} \tag{3}$$

The digram approximation F_2 gives the result

$$\begin{aligned} F_2 &= - \sum_{i,j} p(i,j) \log_2 p(i,j) \\ &= - \sum_{i,j} p(i,j) \log_2 p(i,j) + \sum_i p(i) \log_2 p(i) \\ &= 7.70 - 4.14 = 3.56 \text{ bits per letter.} \end{aligned} \tag{4}$$

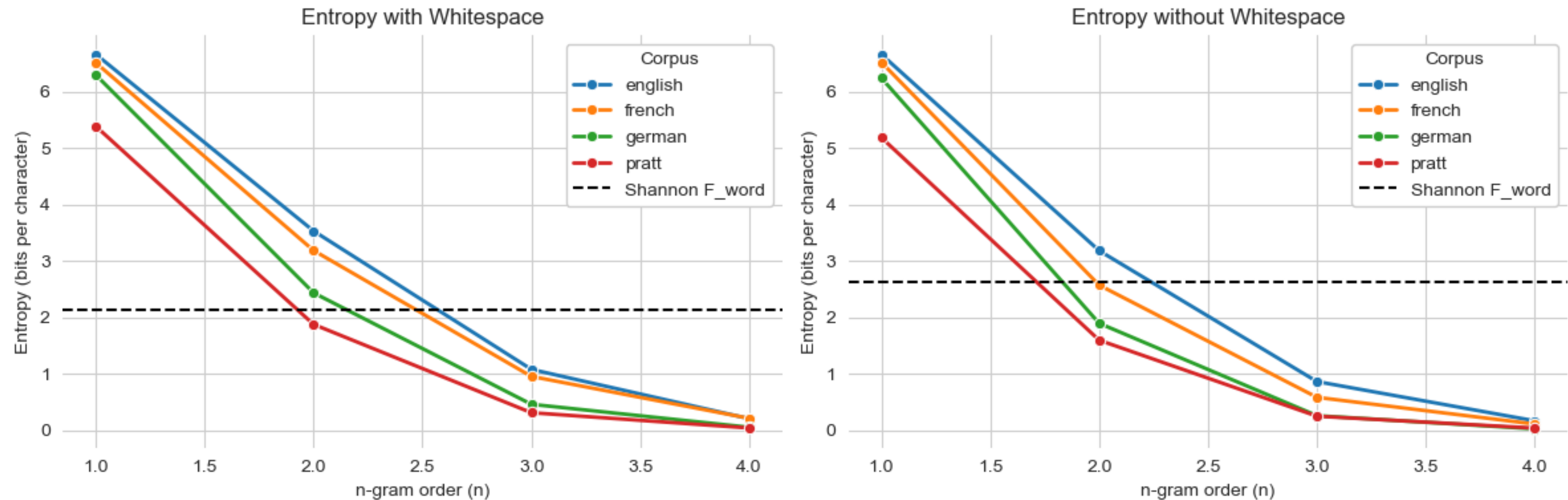
² Fletcher Pratt, "Secret and Urgent," Blue Ribbon Books, 1942.

	F_0	F_1	F_2	F_3	F_{word}
26 letter.....	4.70	4.14	3.56	3.3	2.62
27 letter.....	4.76	4.03	3.32	3.1	2.14

Results from Shannon Paper

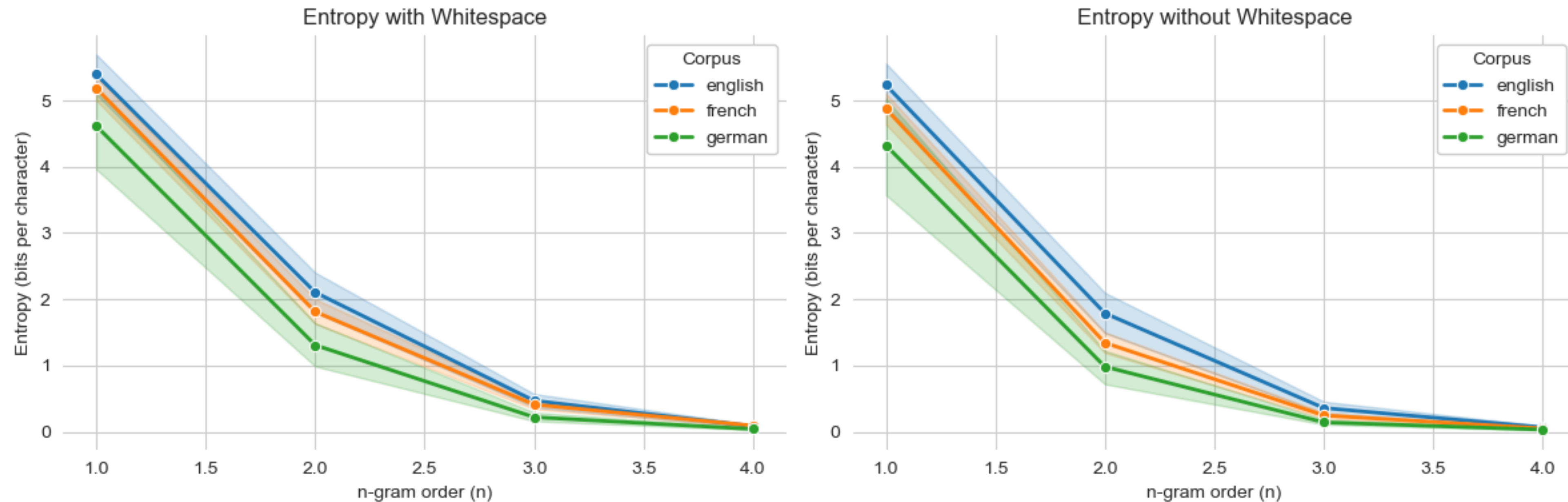
Char-level Entropy Across Languages

Comparison with tabulated values from Shannon

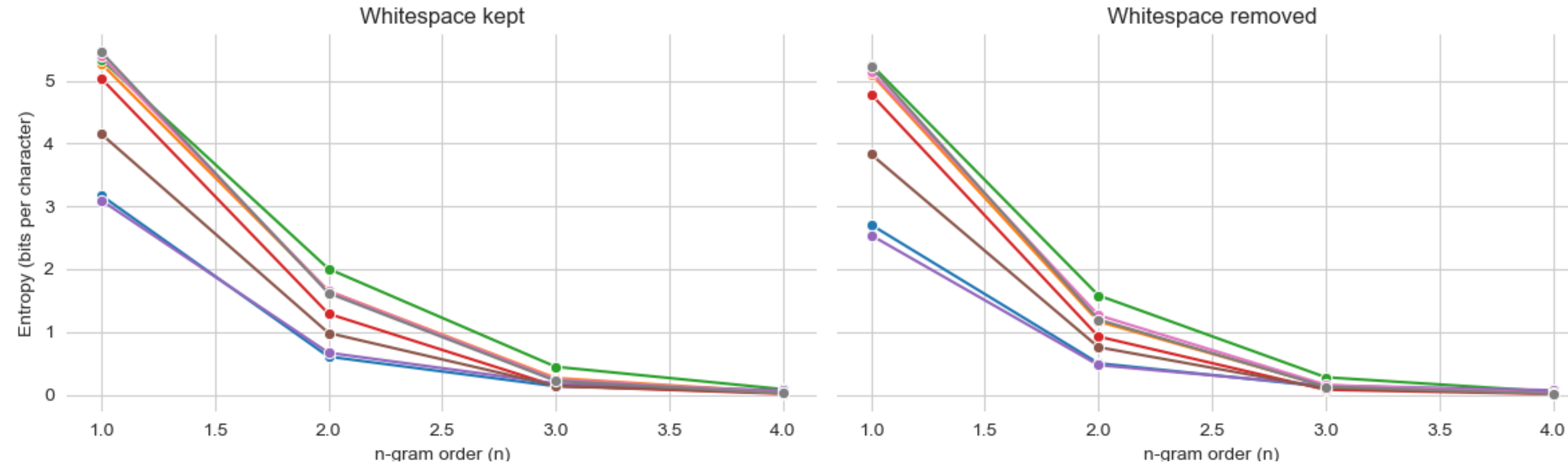


Char-level Entropy Across Languages

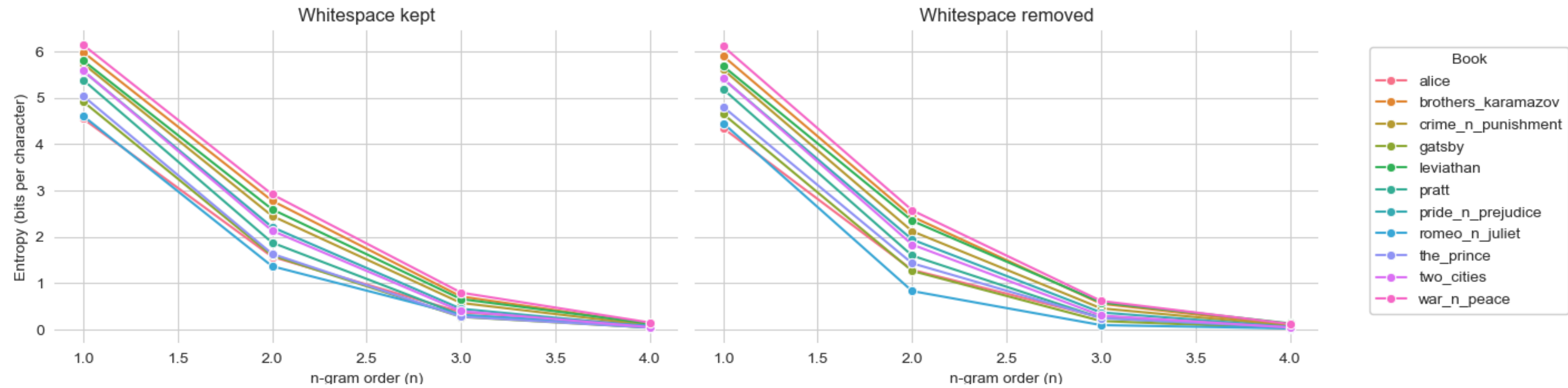
Book level analysis for each language



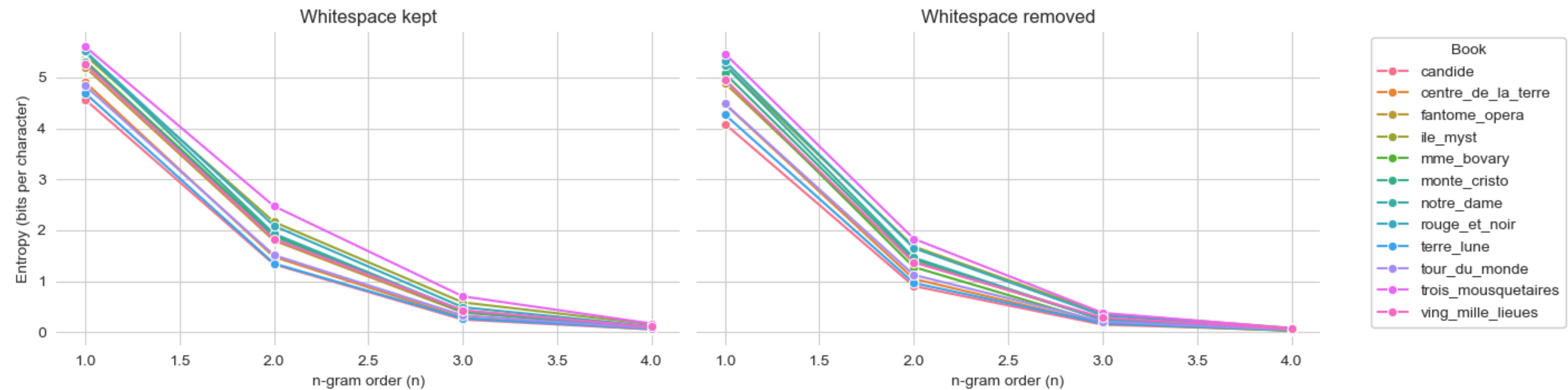
German books



English books



French books



Excursion to Chinese

- Problem = Chinese does not have whitespaces between words
- therefore character-level n-grams are possible, but word-level n-grams are slightly more involved
- luckily there are libraries that do word segmentation by most frequently occurring words -> eg jieba
- 我来到巴黎综合理工学院 = I have come to École Polytechnique in Paris
 - 我 + 来到 + 巴黎 + 综合 + 理工 + 学院
 - 我 + 来到 + 巴黎 + 综合 + 理工学院

Example Model Output

《红楼梦》 Dream of the Red Chamber (18th-century Chinese novel)

----- Model Statistics -----

Order of N-gram Model (N)	: 2
Computation Time (s)	: 0.3320
Number of Unique Contexts:	: 144,941
Total Observed N-grams (Transitions)	: 812,964
Unique (Context, Token) N-grams	: 426,849
Conditional Vocab Size (Next Tokens)	: 4,270

Sample counts:

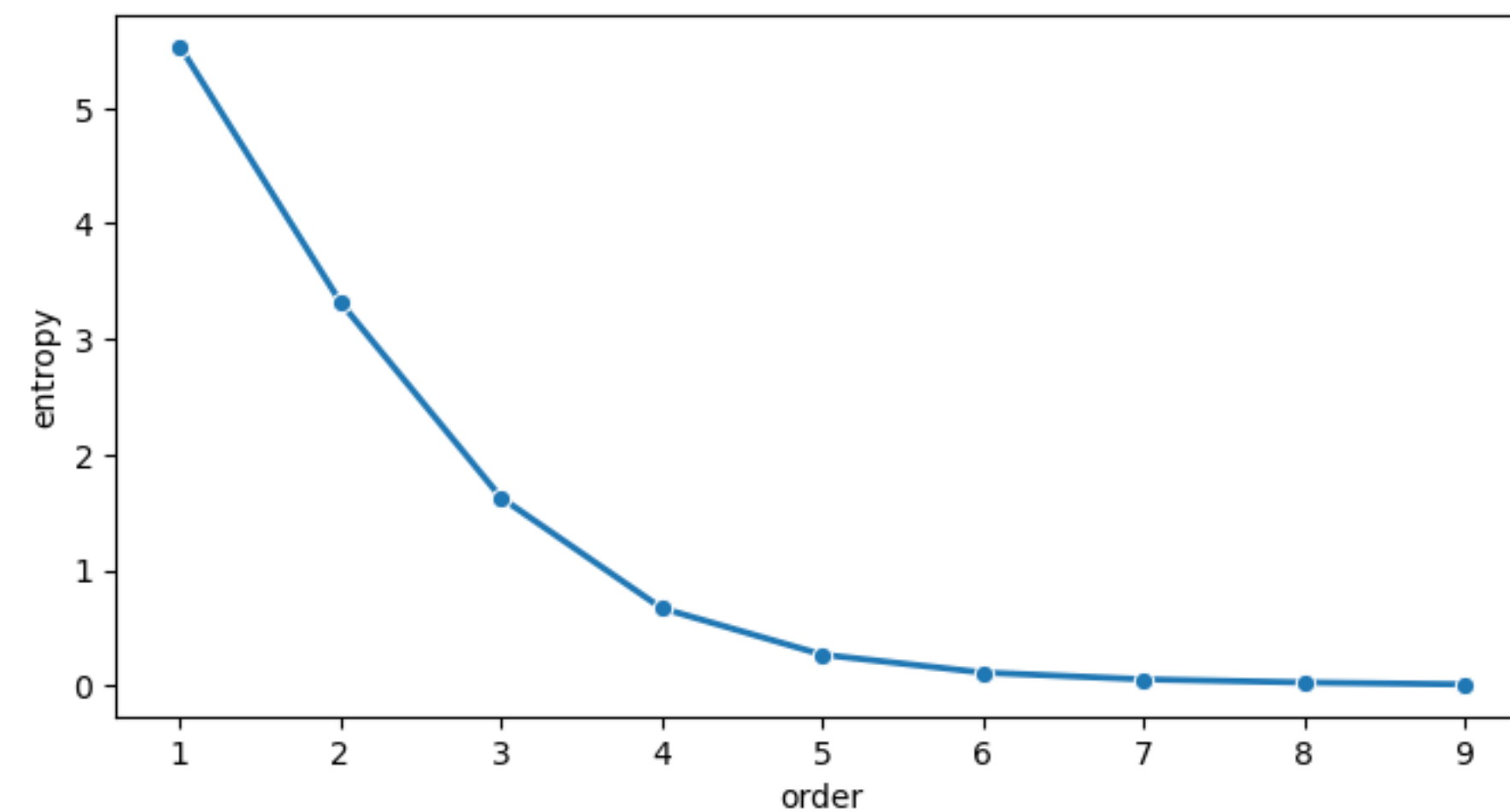
```
[(('紅', '樓'), Counter({'夢': 9})),  
 (('樓', '夢'),  
  Counter({'第': 2, '。': 2, '仙': 1, '稿': 1, '引': 1, '通': 1, '，': 1, '話': 1}))]
```

Most likely token after ('來', '，') :

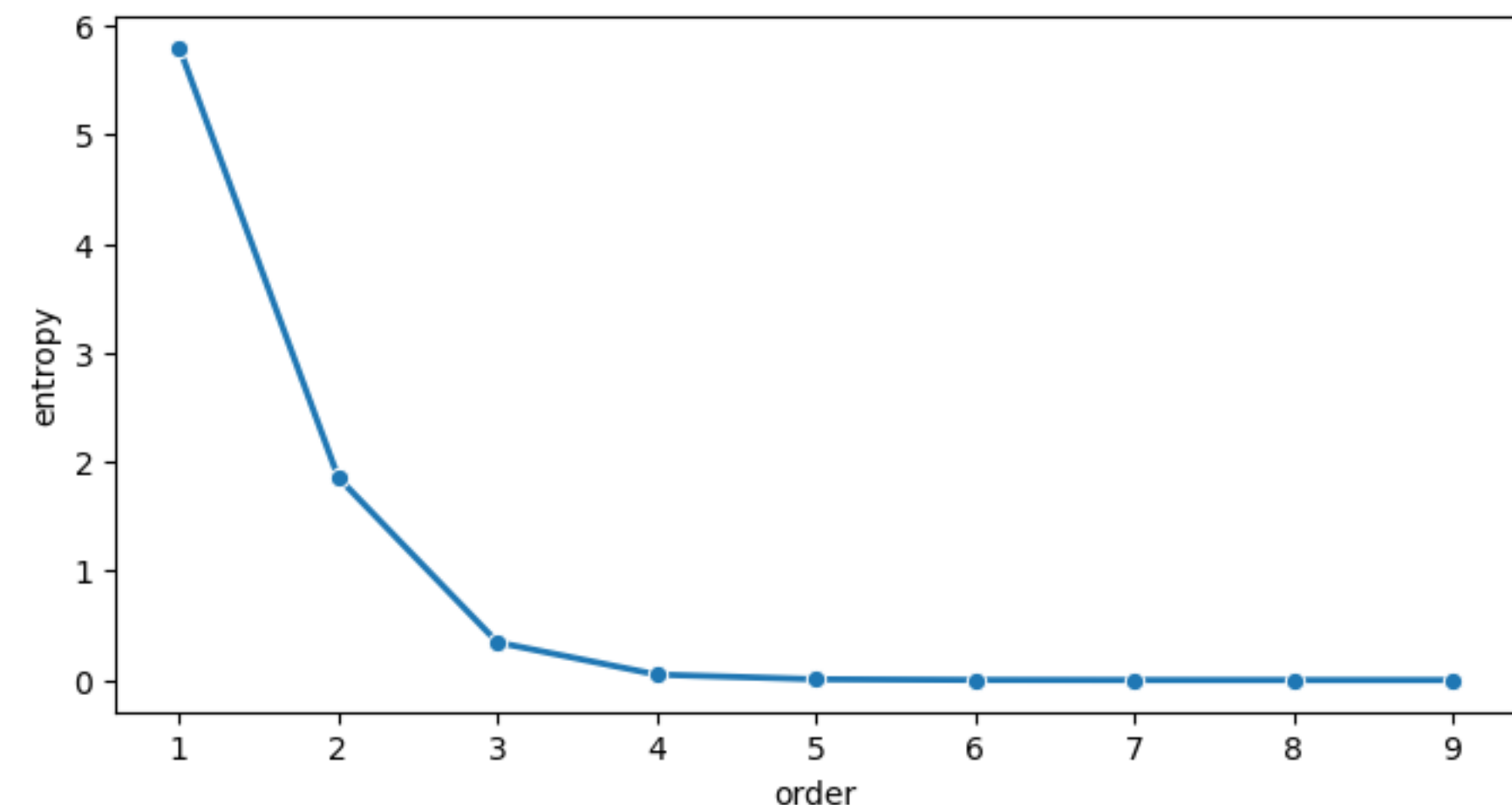
	token	proba
64	我	0.035526
21	說	0.032566
11	只	0.029934
39	便	0.028618
26	—	0.027632
..
331	階	0.000329
332	謝	0.000329
335	淘	0.000329
338	賴	0.000329
539	執	0.000329

Chinese Char- and Word-level Entropy

《红楼梦》 Dream of the Red Chamber (18th-century Chinese novel)



Char-level Entropy in Chinese



Word-level Entropy in Chinese

5) Prediction, Illustration, and Limits

Generated Text

Char-level, Order: 0

nncs oafa s hiw dhwdnaeahtuf haideflfs t v ieiasi hd ali sbnrtslssr kwfhdiia dtyo nosoe h an iieare h ilfly l
oo nd o twxai leeip r rtns astase nneeeiuhniihq tsyharghtltfymnkuet tr ttb t yr a

Char-level, Order: 1

ate ysne as nd chiclangod mo rrer shinido n ar t len s anowoupyot r nd wen tenere ins
terayerureacomokixpond rsestofi santhe ju gherouthe inghad tupo s ghi wh the t tr s herlincaise ieg s rdid
tshensof

Char-level, Order: 2

gforyphould lit isher for to the se of threepeal he cousenchou ke mark there it alselself wed the the sh
sheriethy sing oreme begaid if a cre why quithe crome the ort of ful off to up a lown what wed v

Char-level, Order: 3

o me more queen head if you not at as she word his it nothe said the the can the madedly you cook ture here
thatter dormoutoff anxieuse doubtful zigzag willar and one enought to get minings despeaking ab

Char-level, Order: 4

oung close the play cross tallest explain sigh her own only because of it is was out chin his so herself to
pleasant alices at made and the tarts the said the duchess she queen want to caller and of its t

Generated Text

Word-level, Order: 0

if matters skurried must may a i perfectly with long or shouting said the here loud it but patiently alice the always youve in answer baby of what to till came duchess was dry its above into tone not that howling and it come said me out quite there a

Word-level, Order: 1

pigeon i did with pink eyes very curious as it had never said the cat remarked the dodo in a red rosetree she sat silent for i wonder she ran round the twentieth time of very deep and the while the duck its a thing said the blame on with the

Word-level, Order: 2

out you know and he called the queen left off sneezing by this time youre nothing but outoftheway things had happened lately that alice said nothing perhaps it was said the caterpillar well perhaps your feelings may be different said alice very humbly you had been broken to pieces please then said

Word-level, Order: 3

just at first the two creatures got so close to her ear youre thinking about something my dear and that makes you forget to talk i cant tell you just now what the moral of that is but i shall remember it in a few minutes she heard a little animal she couldnt

Word-level, Order: 4

she was surprised to find quite a large crowd collected round it there was a dispute going on between the executioner the king and the executioner ran wildly up and down looking for it while the rest of the party went back to the table half hoping she might find another key on it

Limitations

N-gram Markov Model

- computational cost grows exponentially
- empirical transition probabilities no smoothing
- entropy estimator $\hat{H}_N \approx \sum_c P(c) H(X \mid c)$
- finite sample bias due to data sparsity at large N