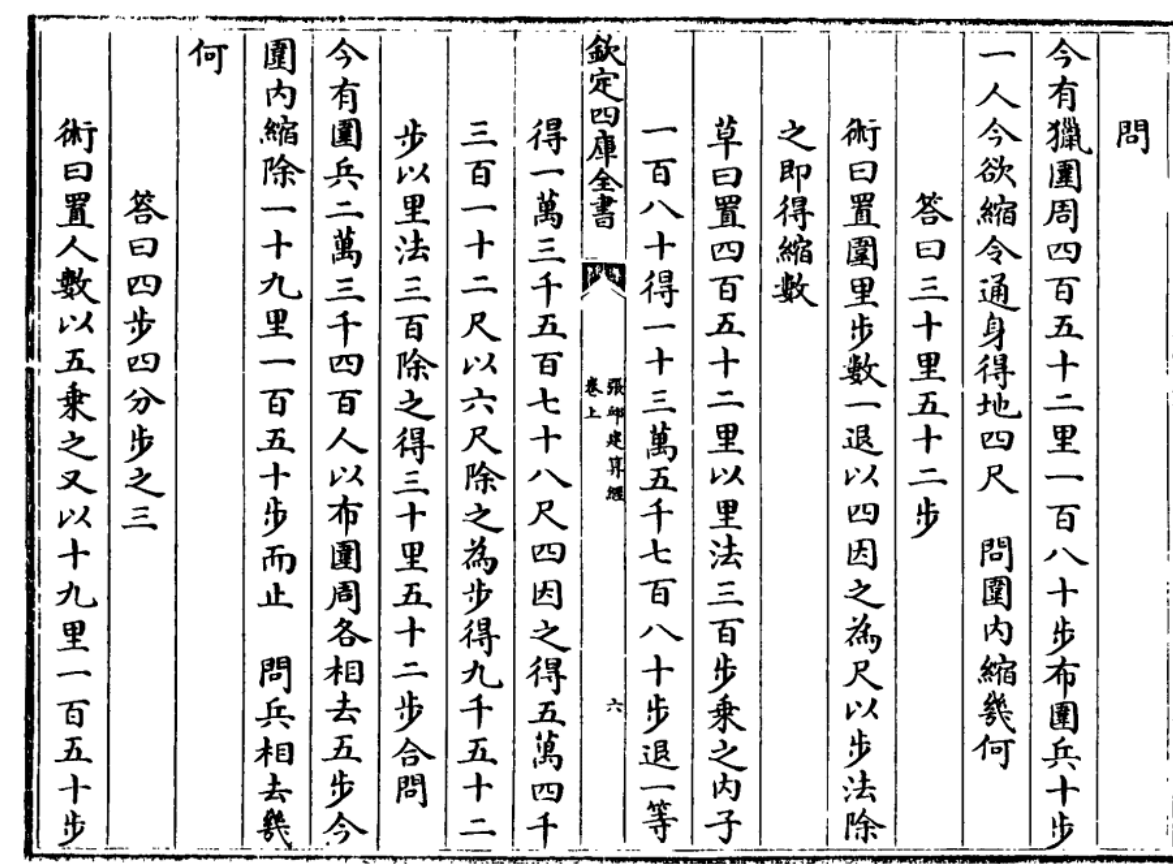
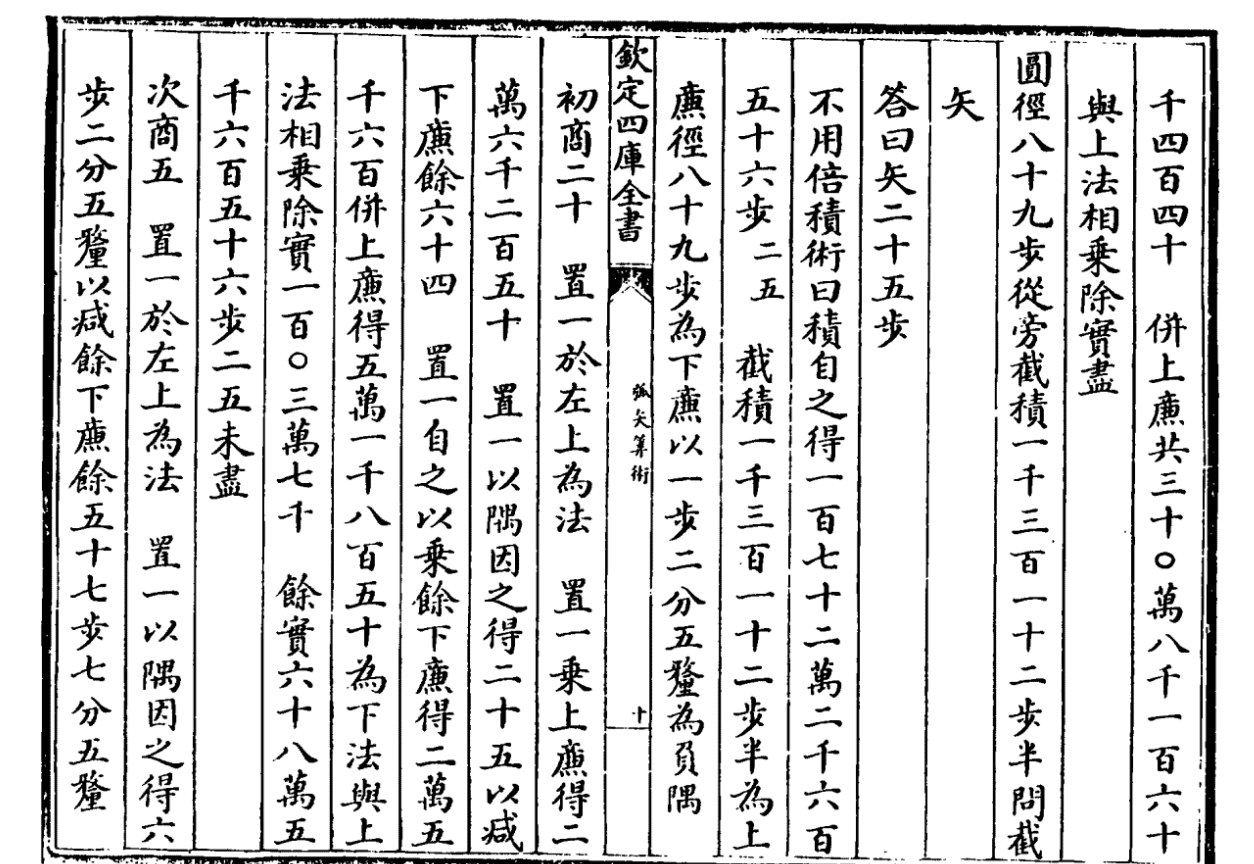


## Pre-Modern Chinese Mathematics: How formulaic was it?



Page from *Siku quanshu* edition of Zhang Qiujian

For some technical genres such as mathematics it is often assumed that their language exhibits **formulaic characteristics**. However, it is notoriously difficult to identify what properties exactly make language “formulaic”. In order to nevertheless quantitatively assess the formulaicity of a corpus, Richard Forsyth has proposed a suite of computational tools titled *formulib*, which operates by “compiling a **formulexicon** from a corpus [...] and then us[es] **coverage** by elements of that formulexicon as an index of the degree to which a text [...] is pervaded by formulaic sequences” (Forsyth, 2021: 33, emphasis ours). In this study, we apply this suite to **pre-modern Chinese mathematical texts** (1st century CE to 18th century). These texts are particularly interesting because they were written **without modern symbolic notation systems**, which could suggest more natural, less formulaic patterns of language use.



Page from *Siku quanshu* edition of Hushi suanshu

## Methodology

### CORPUS

- Use the “Masters Division” (*zibu* 子部) of the *Complete Library of the Four Treasuries* (*Siku quanshu* 四庫全書, compiled 1772-1782).
- Contains different genres: philosophy, arts, and different sciences (including mathematics), ...

### PREPROCESSING

- Remove pre- and postface chapters
- Replace counting rods and numerals with placeholders
- Sampling to obtain balanced corpus, repeat ten times for validation
- After sampling: 11 subcategories with 430 110 tokens each (one Chinese character is one token)

### EXTRACT N-GRAMS

- For each genre, compute *formulexicon*: 80 most frequent *n*-grams,  $3 \leq n \leq 8$
- Exclude *n*-grams made up only of stopwords (= 30 most frequent tokens)

### COMPUTE COVERAGE

- Coverage is the number of tokens of the text that occur in the text as a part of at least one of the *n*-grams in the *formulexicon*, divided by the total number of tokens

## Example: Procedure from Zhang Qiujian's Mathematical Classic

(written ca. 450 CE)

置南北壁高併之得  $N$  半之得  $N$  尺  $N$  寸  
place south north wall height add it obtain NUMERAL halve it obtain NUMERAL *chi* NUMERAL *cun*  
“Place the heights of the southern and norther walls. Add them, obtaining *a*. Halve that, obtaining *b chi* and *c cun*.”

又置長  $N$  尺以廣  $N$  尺因之得  $N$  尺  
furthermore place length NUMERAL *chi* with width NUMERAL *chi* single.digit.multiply it obtain NUMERAL *chi*  
“Furthermore, place the length, *d chi*. Multiply that with the width, *e chi*, obtaining *f chi*.”

又以高  $N$  尺  $N$  寸乘之得  $N$  尺  
furthermore with height NUMERAL *chi* NUMERAL *cun* multiply it obtain NUMERAL foot  
“Furthermore, multiply it with the height, *b chi* and *c cun*, obtaining *i chi*.”

以解法  $N$  尺  $N$  寸  $N$  分除之得  $N$  斛餘  $N$   
with *hu* divisor NUMERAL *chi* NUMERAL *cun* NUMERAL part divide it obtain NUMERAL *hu* remain NUMERAL  
“Divide it by the divisor for *hu*, *j chi*, *k cun* and *l* tenth-*cun*, obtaining *m* with a remainder of *n*.”

## Explanation

All placeholders *a, b, c, ...* are specific numerals in the original

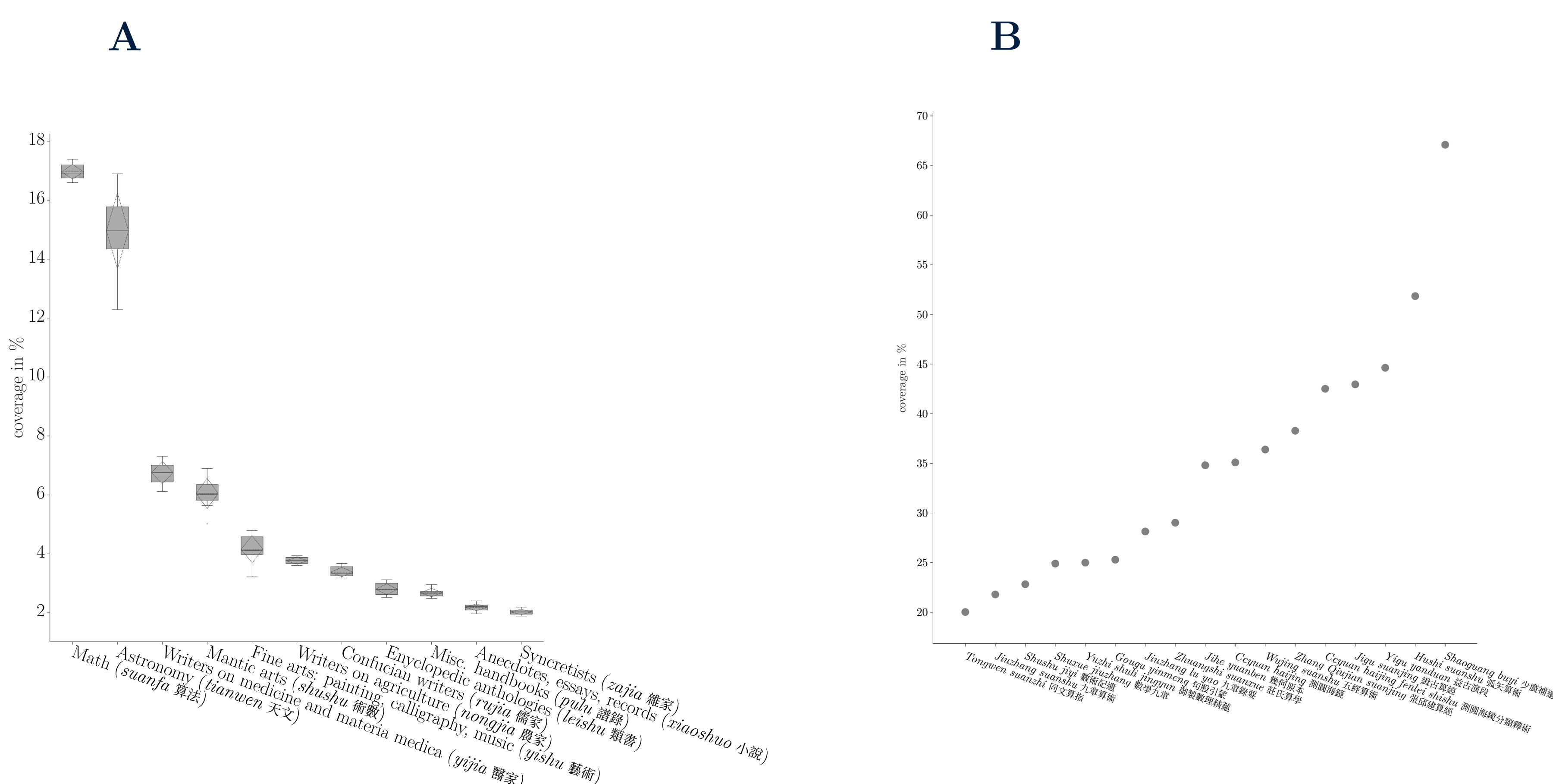
Relevant *n*-grams in the formulexicon:

- Operations: multiplication by single digit (“multiply it, obtaining” 因之得, “multiply it, obtaining *n*” 因之得  $N$ ), regular multiplication (“multiply it, obtaining” 乘之得, “multiply it, obtaining *n*” 乘之得  $N$ ), division (“divide it, obtaining” 除之得, “divide it, obtaining *n* 除之得  $N$ )
- Quantities: “*n chi n*”  $N$ 尺 $N$ , “*n chi n cun*”  $N$ 尺 $N$ 寸, “*n chi n cun n*”  $N$ 尺 $N$ 寸 $N$ , “*n chi n cun n* parts”  $N$ 尺 $N$ 寸 $N$ 分
- Other: “it, obtaining *n chi*” 之得 $N$ 尺, “obtaining *n chi*” 得 $N$ 尺

Coverage calculation: all 30 tokens marked in color are in at least one of the *n*-grams in the formulexicon, so coverage =  $\frac{30}{58} \approx 52\%$

## Results

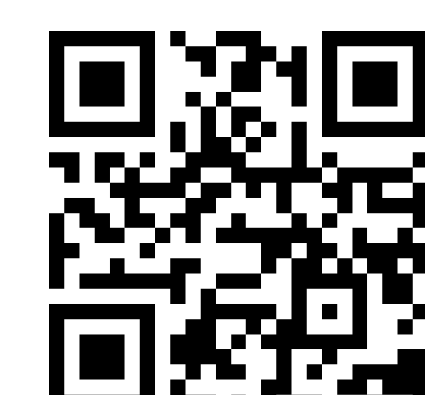
- Coverage per genre shows **mathematical language is more formulaic than other genres**
- Coverage per title (only mathematical texts) shows **high variation in formulaicity**
- Document frequency of *n*-grams in formulexicon shows **low number of formulaic sequences shared by many works**



## Key References

- Forsyth, Richard. 2021. “Cascading collocations: Collocades as correlates of formulaic language”. In *Formulaic language - Theories and Methods*, herausgegeben von Aleksandar Trklja und Lukasz Grabowski. Berlin: Language Science Press.
- Wittern, Christian. 2016. „Kanseki Repository“. CIEAS Research Report 2015 Special issue: 1–80.

Project website:



GitHub:

