

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

Luigi Bambaci¹

¹University of Bologna
Department of Cultural Heritage

12 January 2020

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Plan

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Qohelet: A Digital Scholarly Edition

The parser is part of a PhD project devoted to the preparation of a **digital scholarly edition** of the biblical book of Qohelet

Goals of the edition

- ▶ *Collatio*: digitalization of readings
- ▶ *Constitutio textus*: critical text (**eclectic edition**)
- ▶ Encoding according to the **Text Encoding Initiative**



<https://tei-c.org/>

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Qohelet: A Digital Scholarly Edition

Goals of the parser

- ▶ The parser is designed to encode the readings of Qohelet contained in medieval mss and printed editions as collated by Kennicott
- ▶ The encoding is used for quantitative analysis
- ▶ To identify stemmatic relationships

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Qohelet: A Digital Scholarly Edition

Goals of the parser

- ▶ The parser is designed to **encode** the readings of Qohelet contained in medieval mss and printed editions as collated by Kennicott
- ▶ The encoding is used for **quantitative analysis**
- ▶ To identify **stemmatic relationships**



<http://cophilab.ilc.cnr.it/>



<http://www.ilc.cnr.it/>

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

The book of Qohelet

The book of Qohelet or Ecclesiastes is one of the books of the Hebrew Bible (V-III a.e.v.)

Sources

1. **Direct sources:** witnesses in Hebrew
 - ▶ Qumran scrolls (II-I a.e.v.)
 - ▶ Hebrew medieval manuscripts and printed editions
2. **Indirect sources:** ancient translations (Versions):
 - 2.1 **Primary translations** (Hebrew source):
 - ▶ Greek (Septuaginta, I-II e.v.)
 - ▶ Syriac (Peshitta, II e.v.)
 - ▶ Latin (Vulgata, *Commentarius*, IV e.v.)
 - ▶ Aramaic (Targumim, VII e.v.)
 - 2.2 **Secondary translations** (Greek source):
 - ▶ Ethiopic (IV e.v.)
 - ▶ Armenian (V e.v.)
 - ▶ Syro-hexapla (VII e.v.)
 - ▶ ...

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

Kennicott's work

Optical character
recognition

The Context-free Grammar

The Visitor



VETERIS · TESTAMENTI
EX IMMENSA
MSS. EDITORUMQ. CODICUM CONGERIE HAUSTAE
ET AD SAMAR. TEXTUM, AD VETUSTES VERSIONES,
AD ACCURATIORES SACRAE CRITICAE FONTES AC LEGES
EXAMINATAE
OPERA AC STUDIO
JOHANNIS BERN. DE-ROSSI S. T. D.
ET IN R. PARMENSIS ACAD. LING. OR. PROFESS.

ISAIAS, JEREMIAS, EZECHIEL, XI PROPHETÆ MINORES,
CANTICUM, RUTH, THRENI, ECCLESIASTES, ESTHER.

Veterum librorum fides de Hebraeis voluminibus
examinanda est.

Hersen, Krize, and Lavin.



PARMAE

EX REGIO TYPOGRAPHEO
c12. 1266. LXXXV L

De Rossi, 1788

Medieval tradition: the XVIII c. collations

- ▶ Kennicott and De Rossi gathered thousands of **variants** from more than **1500 witnesses** of the Hebrew Bible
- ▶ Kennicott:
 - 2 volumes, ca. 1800 pages
 - 600 witnesses
 - 1.500.000 pieces of textual information (Barthélemy)
- ▶ The collations contain important **text-critical** and **linguistic** information about the history of the biblical text and Hebrew language in the Middle Ages
- ▶ No extensive collations have been planned since then
- ▶ They are available only in a digitized format (.pdf)

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Qohelet according to Kennicott

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

LIBER ECCLESIASTÆ.

C A P U T I.

11 אִין זוכרן לראשנים וגם לאחרנים שיהיו לא
12 יהיה להם זכרון עמי שיהיו לארצתה: ואני קהלת
13 יהיה מלך על ישראל בירושלם: ונתתי את לבי
14 ללדש ולתור בחכמה על כל אשר נעשה חת
15 השמים: הוא ענין דר נתן אליהם לבני האדם
16 לענות בו: ראיתי את כל המעשים שנעשו תחת
17 השמש והנה הכל רועות רוח: מעות לא
18 יוכל לחקן וחסרון לא יוכל להמנות: דברתי אני
עם בני לאמר אני הנה הגדלתי ויוספתתי חכמה
על כל אשר היה לפני על ירושלם וליבי ראה
19 הרבה חכמה רדע: ואונתו לבי לדעת חכמה
ודעת החלות ושכלות ידעתי שגם זה הוא רעיון
20 רוח: כי ברב חכמה בר כעס ויוסף דעת ויוסף
כסבות:

VARIÆ LECTIONES.

1. דברי ר' קהר' lit. majorit. 4, 109. דברי vox major, et or-
nata; 136, 139 — non major; 1, 2, 3, 4, 21, 57, 67, 82, 89, 93,
99, 100, 110, 118, 128, 130, 141, 144, 231, 237, 239, 270, 280,
ירושלים — דרית — דרד 121. דרד 57, 100, 260; forte 141.
107, 109, 152 — sup. rap. 139. ירודה ירושלים — 76.
31. חביל forte 31 חב' 14. חבל חביל 31 bis 99.
31. חביל — לאיש — לאיש 147. עטל 31 שימול 31.
166, 601.

166, 172, 173, 176, 188, 191, 196, 201, 202, 213, 218, 224, 226, 227, 228, 231, 239, 240, 245, 252, 253, 256, 275, 284, 680, 693, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000.

► Reference text: edition of Ev. **van der Hooght** (1710)

Kennicott's work

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Qohelet according to Kennicott

LIBER ECCLESIASTÆ.

C A P U T I.

11 אִין זוכרן לראשנים וגם לאחרנים שיהיו לא
12 יהיה להם זכרון עמי שיהיו לארצתה: ואני קהלת
13 יהיה מלך על ישראל בירושלם: ונתתי את לבי
14 ללדש ולתור בחכמה על כל אשר נעשה חת
15 השמים: הוא ענין דר נתן אליהם לבני האדם
16 לענות בו: ראיתי את כל המעשים שנעשו תחת
17 השמש והנה הכל רועות רוח: מעות לא
18 יוכל לחקן וחסרון לא יוכל להמנות: דברתי אני
עם בני לאמר אני הנה הגדלתי ויוספתתי חכמה
על כל אשר היה לפני על ירושלם וליבי ראה
19 הרבה חכמה רדע: ואונתו לבי לדעת חכמה
ודעת החלות ושכלות ידיעתי שגם זה הוא רעיון
20 רוח: כי ברב חכמה בר כעס ויוסף דעת ויוסף
כסבות:

VARIÆ LECTIONES.

1. דברי vox major, et orna-
 menta; 136, 139, non major; 1, 2, 3, 14, 31, 57, 67, 82, 89, 93,
 99, 100, 110, 118, 128, 130, 141, 144, 231, 237, 239, 270, 280,
 בירושלים 121. דרד דרד 57, 100, 260; forte 141.
 קדחת – רחית – דרד 107, 109, 152 – sup. raf. 139.
 תורה בירושלים – הבל 31, 14.
 חביל 31. forte חביל 31.
 bis 99.
 הבל הבלין
 עטל 147.
 לאיש – לאדם.
 עטל 147.

166, 172, 173, 176, 188, 191, 196, 201, 202, 213, 218, 224, 226, 227, 228, 231, 239, 240, 245, 252, 253, 259, 275, 284, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000.

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

Kennicott's work

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

- ▶ Reference text: edition of Ev. **van der Hooght** (1710)
- ▶ Critical apparatus: **350** witnesses and ca. **2600** variants

Digitalizing Kennicott's collation

Digitalization can provide useful information for:

► Textual history and philology

- Computing genealogical relationships (**stematology**)
[Hempel, Goshen-Gottstein, Gese, Sacchi, Borbone]
- Studying the phenomenology of the copying process
(scribal habits, genesis of common copying errors...)

► Codicology and paleography

- Classifying mss according to ethno-geographic criteria
(Ashkenazic/Sephardic/Italian mss...)

[Penkower 1988, 2002]

► Linguistics

- *usus scribendi* (*scriptio plena/defectiva*, orthography...)

[Cohen 1986]

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Digitalizing Kennicott's collation

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

In order to permit the computer to extract information from the critical apparatus, data need to be not only **machine readable**, but also fully **machine actionable**

Machine readability is achieved through **OCR** technology (**digitization**)

Machine actionability is ensured by **textual encoding** through **markup languages** (**XML**)

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Towards an automated encoding

- ▶ Manually encoding critical apparatus is expensive, time-consuming and error-prone
- ▶ It is possible to encode **automatically** through Natural Language Processing (**NLP**) tools
- ▶ There are two main approaches:
 1. **Rule-based** systems:
 - ▶ Rules for describing the language are defined by the user
 - ▶ Robust, but the language needs to be as **structured** as possible
 - ▶ Relatively simple
 2. **Machine learning** systems:
 - ▶ Rules for describing the language are derived from the data (machine learning algorithms)
 - ▶ More efficient for non- or semi-structured languages
 - ▶ Complex

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

The language of Kennicott's apparatus

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

- ▶ Well **structured**
- ▶ Textual phenomena are expressed in a **non-redundant** and **unambiguous** way by means of:
 1. Conventional **vocabulary**: **set of finite symbols** (numbers, strings, abbreviations)
 2. Rigorous **syntax**: the **position** of apparatus components conveys information about their function

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Parsing the critical apparatus

An example from Qoh. 1:1

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם. 1.

Words of Qohelet, son of David king in Jerusalem

Kennicott's apparatus

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

76. יהודה בירושלם — sup. ras. 139 — 107, 109, 152 בירושלם

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם. 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

76. יהודה בירושלים — sup. ras. 139 — 107, 109, 152 בירושלים.

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

76. יהודה בירושלים — sup. ras. 139 — 107, 109, 152 בירושלים

variation place (number of chapter and verse)

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

יהודה בירושלים 76. — sup. ras. 139 — 107, 109, 152 בירושלים

readings

דברי קהלת בן דוד מלך בירושלם 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

יהודה בירושלים 76. — sup. ras. 139 — 107, 109, 152 בירושלים

witnesses sigla

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלית — קהלת 121.

דויד 57, 100, 260 ; forte 141.

76. יהודה בירושלים — sup. ras. 139 — 107, 109, 152 בירושלים.

variant description (*primo, nunc, forte* etc.)

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

דברי קהלת בן דוד מלך בירושלם 1.

1. דברי קהלת lit. majorib. 4, 109.

קהלת — קהוית 121.

דויד 57, 100, 260; forte 141.

76. יהודה בירושלים — 139 sup. ras. — 152, 109, 107 בירושלים

separators

Kennicott's work

Optical character
recognition

The Visitor

TEI encoding

lemma (reading from the reference text)

Apparatus components

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

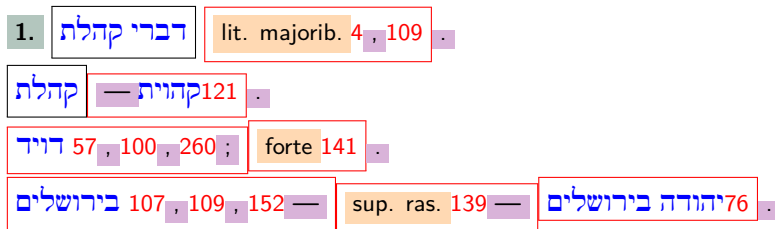
The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

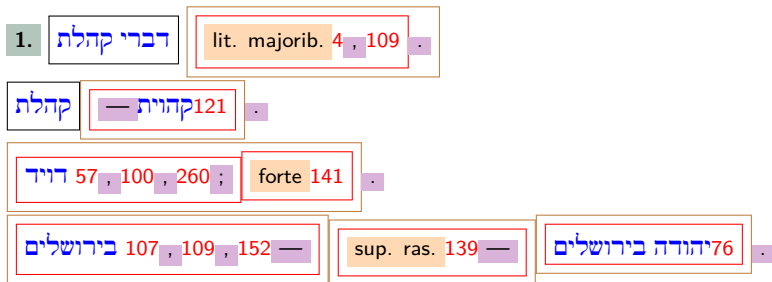
דברי קהלת בן דוד מלך בירושלם. 1.



variant readings diverging from the lemma

Apparatus components

דברי קהלת בן דוד מלך בירושלם. 1.



reading groups

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parsing system

- ▶ It is possible to exploit the nature of structured language characterizing Kennicott's apparatus in order to perform an **automated encoding**
- ▶ The machine is instructed on how to recognize the apparatus components through a **parser**
- ▶ A parser is a software that analyses sequences of strings (**parsing**) according to given rules (**rule-based**)

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parsing system: pipeline

1. The apparatus is processed through **OCR technology** to make it machine readable
2. A **context-free grammar (CFG)** is built to describe the language of the apparatus
3. A **Visitor** is implemented to produce XML code
4. **XSL-T** stylesheets convert XML to TEI

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parsing system: pipeline

1. The apparatus is processed through **OCR technology** to make it machine readable
2. A **context-free grammar (CFG)** is built to describe the language of the apparatus
3. A **Visitor** is implemented to produce XML code
4. **XSL-T** stylesheets convert XML to TEI

ANTLR 4 (*AN*other *T*ool for *L*anguage *R*ecognition)



<https://www.antlr.org/>

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

2. The Context-free grammar

- ▶ The Context-free grammar (CFG) is a formal grammar consisting of a set of top-down, rewriting **rules**
- ▶ The rules of a CFG describe a formal language
- ▶ A **lexer** and a **parser** are implemented in ANTLR4 on the base of the CFG
- ▶ The lexer tokenizes the formal language (**tokenization rules**)
- ▶ The parser checks its syntax (**parser rules**) creating an **Abstract Syntactic Tree (AST)**

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor
TEI encoding

Stemmatic analysis

2. The Context-free grammar

The readings

דברי קהלת lit. majorib. 4, 109.

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor
TEI encoding

Stemmatic analysis

2. The Context-free grammar

The readings

קהלת lit. majorib. 4, 109.

```
7 grammar kennicottCFG;
8 all: listApp+;
9 listApp: loc app+ ;
10 app: lem? rdgGrp+ closeApp;
11 lem: w+ lemSep;
12 rdgGrp: (rdg+) rdgGrpSep?;
13 rdg: (w+)? (term)? wits rdgSep?;
14 w: HEBW;
15 loc: verse closeLoc;
16 term: MAN_DESC;
17 wits : wit+;
18 sigl: NUM;
19 wit: sigl com?;
20 lemSep: VAR_SEP;
21 rdgGrpSep: VAR_SEP ;
22 com: COMMA;
23 rdgSep: COMMA | SEMICOLON;
24 numSign: NUMEROSIGN ;
25 verse: num+ ;
26 num: NUM;
27 closeMainApp: NEWLINE;
28 closeApp: END TAB | END NEWLINE ;
29 closeLoc: END;
30
31 MAN_DESC: 'forte' | 'sup. ras.' |
32 'lit. majorib.' ;
33 ALPHA_SEQ : [a-zA-Z]+;
34 NUM : [0-9]+('.'[0-9]+)?;
35 HEBW : [\u0590-\u05ff*]+;
36 NEWLINE: ('\n');
```

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

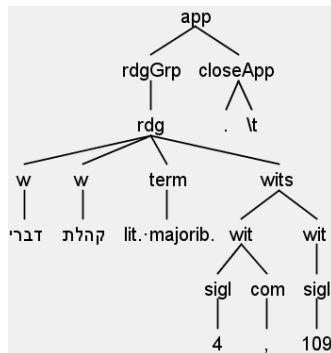
Stemmatic analysis

2. The Context-free grammar

The readings

דברי קהלת lit. majorib. 4, 109.

```
7 grammar kennicottCFG;
8 all: listApp+;
9 listApp: loc app+ ;
10 app: lem? rdgGrp+ closeApp;
11 lem: w+ lemSep;
12 rdgGrp: (rdg+) rdgGrpSep?;
13 rdg: (w+)? (term)? wits rdgSep?;
14 w: HEBW;
15 loc: verse closeLoc;
16 term: MAN_DESC;
17 wits : wit+;
18 sigl: NUM;
19 wit: sigl com?;
20 lemSep: VAR_SEP;
21 rdgGrpSep: VAR_SEP ;
22 com: COMMA;
23 rdgSep: COMMA | SEMICOLON;
24 numSign: NUMEROSIGN ;
25 verse: num+ ;
26 num: NUM;
27 closeMainApp: NEWLINE;
28 closeApp: END TAB | END NEWLINE ;
29 closeLoc: END;
30
31 MAN_DESC: 'forte' | 'sup. ras.' |
32 'lit. majorib.' ;
33 ALPHA_SEQ : [a-zA-Z]+;
34 NUM : [0-9]+('.'[0-9]+)?;
35 HEBW : [\u0590-\u05ff*]+;
36 NEWLINE: ('\n');
```



Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

2. The Context-free grammar

The reading groups

76. יהודה בירושלים — sup. ras. 139 — 107, 109, 152 בירושלים

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

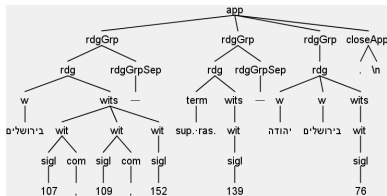
Stemmatic analysis

2. The Context-free grammar

The reading groups

76. יהודה בירושלים — sup. ras. 139 — 107, 109, 152 בירושלים

```
7  grammar kennicottCFG;
8  all: listApp+;
9  listApp: loc app+ ;
10 app: lem? rdgGrp+ closeApp;
11 lem: w+ lemSep;
12 rdgGrp: (rdg+ ) rdgGrpSep?;
13 rdg: (w+)? (term)? wits rdgSep?;
14 w: HEBW;
15 loc: verse closeLoc;
16 term: MAN_DESC;
17 wits : wit+;
18 sigl: NUM;
19 wit: sigl com?;
20 lemSep: VAR_SEP;
21 rdgGrpSep: VAR_SEP ;
22 com: COMMA;
23 rdgSep: COMMA | SEMICOLON;
24 numSign: NUMEROSIGN ;
25 verse: num+ ;
26 num: NUM;
27 closeMainApp: NEWLINE;
28 closeApp: END TAB | END NEWLINE ;
29 closeLoc: END;
30
31 MAN_DESC: 'forte' | 'sup. ras.' |
32 'lit. majorib.' ;
33 ALPHA_SEQ : [a-zA-Z]+;
34 NUM : [0-9]+('.'[0-9])?;
35 HEBW : [\u0590-\u05ff]+;
36 NEWLINE: ('\n');
```



Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

3. The Visitor

- ▶ The Visitor is a tree-walking mechanism
- ▶ In this implementation, the Visitor traverses the tree and slavishly translates the CFG rules into XML nodes
- ▶ It is a **general-purpose** exporter: once implemented, it can be applied to every CFG without further customization

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

3. The Visitor

- ▶ The Visitor is a tree-walking mechanism
- ▶ In this implementation, the Visitor traverses the tree and slavishly translates the CFG rules into XML nodes
- ▶ It is a **general-purpose** exporter: once implemented, it can be applied to every CFG without further customization

<https://cophilab.ilc.cnr.it/parseForge/>
(Courtesy of F. Boschetti)

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

3. The Visitor: from the AST to XML

Introduction

Qohelet
Kennicott's work

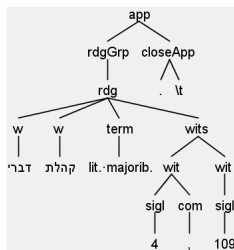
Digitalization

Optical character
recognition
The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis



```
1 <listapp>
2   <loc>
3     <verse>
4       <num>1</num>
5     </verse>
6     <closeLoc>.</closeLoc>
7   </loc>
8   <app>
9     <rdgGrp>
10      <rdg>
11        <w>דבר</w>
12        <w>קהלת</w>
13        <term>lit. majorib.</term>
14        <wit>
15          <sigl>4</sigl>
16          <com>,</com>
17        </wit>
18        <wit>
19          <sigl>109</sigl>
20        </wit>
21      </rdg>
22    </rdgGrp>
23  </app>
24 </listapp>
25 ...
```

4. XSL-T: from XML to TEI

Location referenced method

```
1 <listapp>
2   <loc>
3     <verse>
4       <num>1</num>
5     </verse>
6     <closeLoc>.</closeLoc>
7   </loc>
8   <app>
9     <rdgGrp>
10      <rdg>
11        <w>דברי</w>
12        <w>קהלת</w>
13        <term>lit. majorib.</term>
14        <wit>
15          <sigl>4</sigl>
16          <com>,</com>
17        </wit>
18        <wit>
19          <sigl>109</sigl>
20        </wit>
21      </rdg>
22    </rdgGrp>
23  </app>
24 </listapp>
25 ...
```



```
1 <listapp>
2   <app loc="1 1">
3     <lem>
4       <w>דברי</w>
5       <w>קהלת</w>
6     </lem>
7     <rdgGrp>
8       <rdg wit="#K4 #K109">
9         <term>lit. majorib.</term>
10      </rdg>
11    </rdgGrp>
12  </app>
13  <app loc="1 1">
14    <lem>
15      <w>קהלת</w>
16    </lem>
17    <rdgGrp>
18      <rdg wit="#K121">
19        <w>קהלית</w>
20      </rdg>
21    </rdgGrp>
22  </app>
23  ...
```

A rule-based parser

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Results

1. 100% of accuracy for Qohelet

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

A rule-based parser

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Results

1. 100% of accuracy for Qohelet
2. 2617 variants correctly parsed and encoded

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parser

Results

1. 100% of accuracy for Qohelet
2. 2617 variants correctly parsed and encoded
3. Rut (1025 variants) and Song of Songs (1238): few syntactic errors (99% accuracy)

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parser

Results

1. 100% of accuracy for Qohelet
2. 2617 variants correctly parsed and encoded
3. Rut (1025 variants) and Song of Songs (1238): few syntactic errors (99% accuracy)
4. The CFG is extended in order to include unseen textual phenomena (100% accuracy)

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

A rule-based parser

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Evaluation

- ▶ The system is robust
- ▶ Other biblical books of Kennicott's collection can be encoded

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

A rule-based parser

Evaluation

- ▶ The system is robust
- ▶ Other biblical books of Kennicott's collection can be encoded
- ▶ A parser approach is **faster** than manual encoding
- ▶ It is possible to export data in **different formats** (XML, HTML, relational databases, L^AT_EX etc.)
- ▶ It is **less error prone**: the CFG provides a **spell-checking system**, useful for detecting
 1. errors generated after OCR (transcriptional errors)
 2. inconsistencies in the printed source
- ▶ Tighter control on **semantic errors** when exporting to XML-TEI

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Treatment of residuals

- ▶ Kennicott's apparatus minimizes **natural language**
- ▶ Encoding by means of **concise** and **standard annotations**

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Treatment of residuals

- ▶ Kennicott's apparatus minimizes **natural language**
- ▶ Encoding by means of **concise** and **standard annotations**
- ▶ Exceptions (residuals):
 - 7:29 "Incipit cap. 8 a voce חכמת, medio commatis..."
 - 11:9 "ביאך" — post hanc vocem sequitur Psal. 102..."

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Stemmatic analysis

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Stemmatic analysis

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

- Encoding can be used for quantitative analysis (e. g. **computer-assisted stemmatology**)

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Stemmatic analysis

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

- ▶ Encoding can be used for quantitative analysis (e. g. **computer-assisted stemmatology**)
- ▶ Few attempts to study the medieval tradition of the Hebrew Bible according to (**neo-**)**Lachmannian** criteria
- ▶ Mss are treated singularly, not as groups or families
- ▶ A **stemmatic classification** of medieval mss is still missing

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

From XML-TEI to data matrix

In order to permit the computer to compute genealogical relationships, it is necessary to transform philological data into a numerical format (**data matrix**)

```
1 <listapp>
2   <app loc="1 1">
3     <lem>
4       <w>דברי</w>
5       <w>קהלת</w>
6     </lem>
7     <rdgGrp>
8       <rdg wit="#K4 #K109">
9         <term>lit. majorib.</term>
10      </rdg>
11    </rdgGrp>
12  </app>
13  <app loc="1 1">
14    <lem>
15      <w>קהלת</w>
16    </lem>
17    <rdgGrp>
18      <rdg wit="#K121">
19        <w>קהלית</w>
20      </rdg>
21    </rdgGrp>
22  </app>
23  ...
```

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Data matrix: regularization

Elimination of

- ▶ Partially collated witnesses
- ▶ Accidentals (e. g. orthographic variants)
- ▶ *Marginalia*, dubious and second hand variants

Total: 119 witnesses and 371 variants

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Phylogenetic analysis

- ▶ The data matrix is analysed by phylogenetic algorithms, which produce tree-like graphs (**phylograms**) representing paths of textual evolution
- ▶ Many phylogenetic algorithms: **Maximum Parsimony**
- ▶ Implemented in **PAUP** software

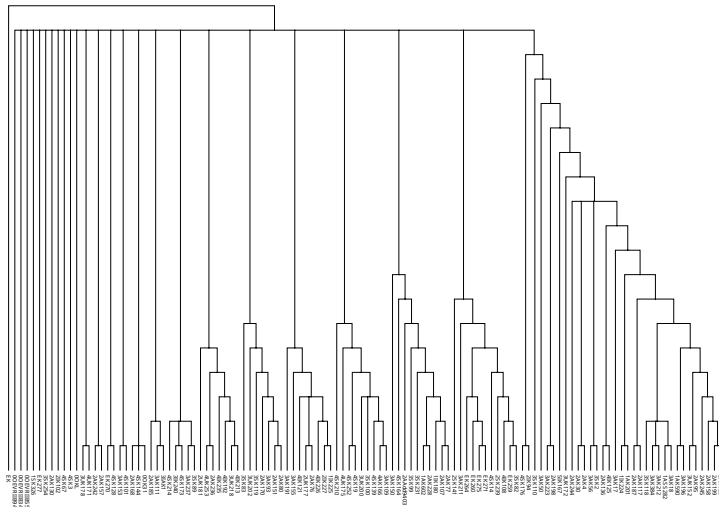
Parsimony Analysis Using PAUP



<http://paup.phylosolutions.com/>

Strict consensus tree

Consensus of more than 25.000 equally parsimonious trees



A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Strict consensus tree

Consensus of more than 25.000 equally parsimonious trees

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

Kennicott's work

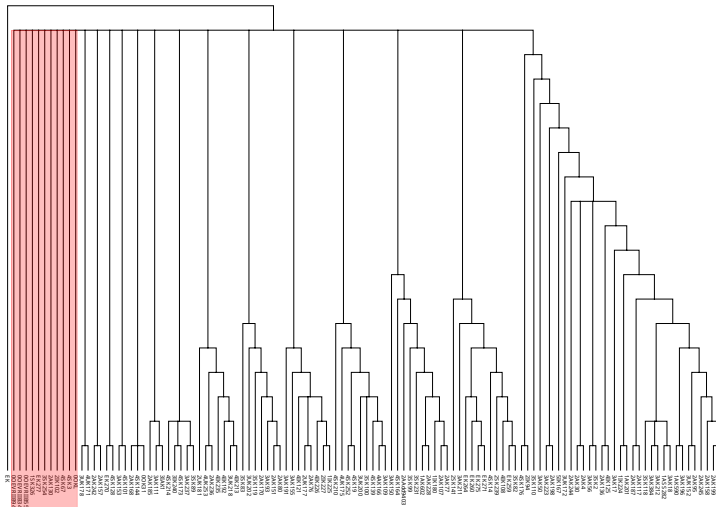
Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis



Strict consensus tree

Consensus of more than 25.000 equally parsimonious trees

A Rule-Based Parsing System for Encoding Kennicott's Collation of the Hebrew Bible

Kennicott's work

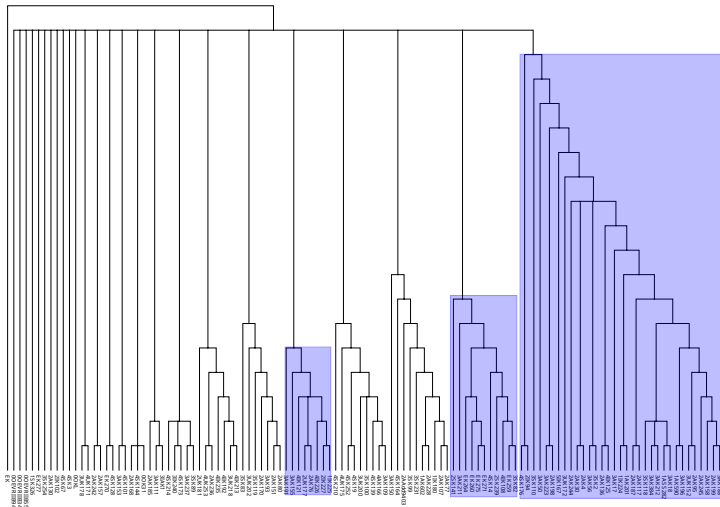
Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis



Phylogenetic methods and biblical mss

- ▶ Groups are identified according to **ancestral variants**
- ▶ Some groups have **characteristic variants**
- ▶ Corroborated by **external criteria**

A Rule-Based
Parsing System
for Encoding
Kennicott's
Collation
of the Hebrew
Bible

Introduction

Qohelet

Kennicott's work

Digitalization

Optical character
recognition

The Context-free Grammar

The Visitor

TEI encoding

Stemmatic analysis

Phylogenetic methods and biblical mss

- ▶ Groups are identified according to **ancestral variants**
- ▶ Some groups have **characteristic variants**
- ▶ Corroborated by **external criteria**

Limits

- ▶ Important mss are missing
- ▶ Variants concern **consonantal text** only
- ▶ Most are **polygenetic**
- ▶ **Contamination**

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Phylogenetic methods and biblical mss

- ▶ Groups are identified according to **ancestral variants**
- ▶ Some groups have **characteristic variants**
- ▶ Corroborated by **external criteria**

Limits

- ▶ Important mss are missing
- ▶ Variants concern **consonantal text** only
- ▶ Most are **polygenetic**
- ▶ **Contamination**

Future works

- ▶ More data:
 - ▶ more mss
 - ▶ punctuation, Massora, para-textual variants
- ▶ Variants can be **weighted** and **ordered**

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis

Bibliographical references

- Barthélemy, D. (1992). Les manuscrits médiévaux et le texte tiberien classique. In *Critique textuelle de l'Ancien Testament*, 3. *Ezéchiel, Daniel et les 12 Prophètes*, Volume 3 of *Orbis Biblicus et Orientalis*, pp. xix–xcvi. Fribourg/Göttingen: Éditions Universitaires/Vandenhoeck & Ruprecht.
- Borbone, P. G. (1990). *Il libro del profeta Osea. Edizione critica del testo ebraico*. Torino: Zamorani.
- Cohen, M. (1973). מנבשי כתיב במצפה מסורה עתיקים ומשמעם לתולדות נוסח המקרא המקובל. Ph. D. thesis, Hebrew University, Jerusalem. [Unpublished].
- Cohen, M. (1979). המקרא ואנחנו. In S. Uriel (Ed.), *המקרא באור*, Volume 1, pp. 42–69. Tel Aviv: דביר / Dvir.
- Cohen, M. (1980). עיוני מקרא ופרשנות. קווי יסוד לדמותו העיצורית של הטקסט בכתבי יד מקראיים מימי הביניים. / *Studies in Bible and Exegesis* 1, 123–182.
- Cohen, M. (1981). לדמותם הקונסוננטית של דפוסי המקרא הראשונים: המהדורה הראשונה של התנ"ך השלם - דפוס 1488 / The Consonantal Character of the First Rabbinic Printings: the *Editio Princeps* of the Entire Bible Soncino 1488. In *ספר השנה של אוניברסיטת בר-אילן*, Volume XVIII-XIX, pp. 47–67. Ramat Gan: University Press.
- Cohen, M. (1986). מהו 'נוסח המסורה' ומה היקף אחיזתו בתולדות המסירה של ימיה' / The 'Masoretic Text' and the Extent of Its Influence on the Transmission of the Biblical Text in the Middle Ages. In S. Uriel (Ed.), *עיוני מקרא ופרשנות* / *Studies in Bible and Exegesis*, Volume 2, pp. 229–256. Ramat Gan: Bar Ilan University Press.
- Gese, H. (1957). Die hebräischen Bibelhandschriften zum Dodekapropheten nach der Variantensammlung des Kennicott. *Zeitschrift für die Alttestamentliche Wissenschaft* 69(1–4), 55.
- Goshen-Gottstein, M. H. (1954). Die Jesaiah-Rolle und das Problem der hebräischen Bibelhandschriften. *Biblica* 35(4), 429–442.
- Hempel, J. (1930). Chronik. *Zeitschrift für die Alttestamentliche Wissenschaft* 48, 187–206.
- Hempel, J. (1934). Innermasoretische Bestätigungen des Samaritanus. *Zeitschrift für die Alttestamentliche Wissenschaft* 52(1), 254–274.
- Penkower, J. S. (1982). העקב בן חיים וצמיחת מהדורת המקראית הגדולה. / *Jacob Ben-Hayyim and the Rise of the Biblia Rabbinica*. Ph. D. thesis, Hebrew University, Jerusalem. [Unpublished].
- Penkower, J. S. (1988). כתב-יד ירושלמי של התורה מן המאה העשירית שהניחו מישאל בן עוזיאל (כתב-יד 3ק). / A Tenth-century Pentateuchal MS from Jerusalem (MS C3), Corrected by Mishael ben Uzziel. *תרביץ* / *Tarbiz* 58(1), 49–74.
- Penkower, J. S. (2002). A Sheet of Parchment from a 10th or 11th Century Torah Scroll: Determining Its Type Among Four Traditions (Oriental, Sefardi, Ashkenazi, Yemenite). *Textus* 21(1), 235–264.
- Sacchi, P. (1973). Analisi quantitativa della tradizione medievale del testo ebraico della Bibbia secondo le collazioni del De Rossi. *Oriens Antiquus* 12, 1–13.
- Wevers, J. W. (1948). A Study in the Hebrew Variants in the Books of Kings. *Zeitschrift für die Alttestamentliche Wissenschaft* 61(1), 43.

Introduction

Qohelet
Kennicott's work

Digitalization

Optical character
recognition
The Context-free Grammar
The Visitor
TEI encoding

Stemmatic analysis