# Dynamic universals in the linguistic marking of location

## extracted from parallel texts

*Michael Cysouw*

Philipps Universität Marburg

# Introducing the massively Parallel Text Corpus

# Worldwide 'survey' data

- **Massively parallel texts**
  - ▸ Same text available in many languages (i.e. translations !)
  - ▸ Contextually situated comparable expressions

- **Including lesser-described languages**
  - ▸ Bible
  - ▸ Universal Declaration of Human Rights
  - ▸ Pamphlets of Jehova's Witnesses

# Bible corpus

- http://paralleltext.info
  - 1169 translations online (soon 1850+)
  - 906 different ISO-639/3 codes (soon 1400+)
  - All data in a private GitHub repo (ask me!)
  - Old Testament 23K verses, New Testament 8K verses, Apocrypha 6K verses
  - All texts cleaned, aligned, normalised, punctuation separated
  - 4 GB raw text files

- Today data from
  - 1556 New Testament translations
  - 1163 different ISO-639/3 codes

# **Single-word comparative linguistics**

The case of 'Jerusalem'

# Angaataha

## (ISO 639-5 agm, spoken in Papua New Guinea)

- jerusaremɨhanda
- jerusaremɨhandaahapɨ
- jerusaremɨhandɨ
- jerusaremɨhandaahiyai
- jerusaremɨhandaahɨ
- jerusaremɨhandaahapɨhiyaunɨ
- jerusaremɨhandaahiyaisangi
- jerusaremɨhandaahapɨhiya
- jerusaremɨhandaahɨraapɨ
- jerusaremɨhandaahɨhɨ
- jerusaremɨhandaahɨhe
- jerusaremɨhandamɨ
- jerusaremɨmanda

- jerusaremɨhandapɨ
- jerusaremɨndɨ
- jerusaremɨhandaahapɨto
- jerusaremɨhandaahapaahɨhɨ
- jerusaremɨhandi
- jerusaremɨmandaahapɨ
- jerusaremɨhandaahunɨ
- jerusaremɨhandaahapunɨ
- jerusaremɨhandaahiya
- jerusaremɨhandamɨhinɨ
- jerusaremɨhandaahapɨhiyaatihɨ
- jerusaremɨhandaahapɨhiyaate
- jerusaremɨhandaahiyaunɨ

# Amharic

## (ISO 639-3 amh, spoken in Ethiopia)

- ኢየሩሳሌም
- በኢየሩሳሌም
- ከኢየሩሳሌም
- ኢየሩሳሌምም
- በኢየሩሳሌምም
- ኢየሩሳሌምን
- ከኢየሩሳሌምም

- የኢየሩሳሌም
- ለኢየሩሳሌም
- ለኢየሩሳሌምም
- የኢየሩሳሌምንም
- የኢየሩሳሌምም

# Amharic

## (ISO 639-3 amh, spoken in Ethiopia)

- ኢየሩሳሌም

- በኢየሩሳሌም

- ከኢየሩሳሌም

- ኢየሩሳሌምም

- በኢየሩሳሌምም

- ኢየሩሳሌምን

- ከኢየሩሳሌምም

- የኢየሩሳሌም

- ለኢየሩሳሌም

- ለኢየሩሳሌምም

- የኢየሩሳሌምንም

- የኢየሩሳሌምም

# Low hanging fruit …

- Semantic structure of locative case

- Complexity of locative case

- Word order patterns

- Phylogenetic correlation between word order and complexity

- Reconstruction of language families based on functional differences

# Low hanging fruit …

- **Semantic structure of locative case**

- Complexity of locative case

- Word order patterns

- Phylogenetic correlation between word order and complexity

- Reconstruction of language families based on functional differences

# Semantic map

- When two contexts recurrently show the same form in language after language, then these contexts have a similar meaning

- Use a low-dimensional reduction of the similarity space to compare languages

'the boy Jesus remained in Jerusalem' (Luke 2:43)

'he was traveling … and making his way to Jerusalem' (Luke 13:22)

'there were Pharisees …, who had come from … Jerusalem ' (Luke 5:17)

Essive ('in')

Inhabitants

Inlative ('into')

Allative ('to')

Apudlative ('near')

Ablative ('from')

Argument

# eng-x-bible-catholic.txt

Jerusalem

Levels drawn at 41-46-51%

# yle-x-bible.txt



Njedusalem

Levels drawn at 41-46-51%

# tur-x-bible-2009.txt



Legend:
- Yeruşalim'e
- Yeruşalim'de
- Yeruşalim'den
- Yeruşalim
- Yeruşalim'in
- other

Levels drawn at 41-46-51%

16

# are-x-bible.txt



Jerusalem-urna
Jerusalem-ala
Jerusalem-arinya
Jerusalem-anga
Jerusalem
other

Levels drawn at 41-46-51%

17

# gug-x-bible.txt



Legend:
- ○ Jerusalénpe
- △ Jerusalén
- + Jerusaléngui
- × Jerusalengua
- ◇ Jerusalenguáva

Levels drawn at 41-46-51%

18

# Low hanging fruit ...

- Semantic structure of locative case

- **Complexity of locative case**

- Word order patterns

- Phylogenetic correlation between word order and complexity

- Reconstruction of language families based on functional differences

**Histogram of entropy**

**Small communities
ignoring languages with zero entro**

**Large communities
ignoring languages with zero entro**

Entropy (square root)

Entropy (square root)

Speaker community size (log10)

Speaker community size (log10)

# Low hanging fruit …

- Semantic structure of locative case

- Complexity of locative case

- **Word order patterns**

- Phylogenetic correlation between word order and complexity

- Reconstruction of language families based on functional differences

**English ( eng-x-bible-catholic.txt )**
**(wordforms)**



Jerusalem

Levels drawn at 41-46-51%

**English ( eng-x-bible-catholic.txt )**
**(before, binding: 18.59 )**



to
in
at
from

Levels drawn at 25%

**English ( eng-x-bible-catholic.txt )**
**(after, binding: 2.33 )**



to
and
the
from

Levels drawn at 25%

**Hindi ( hin-x-bible-easy.txt )**
**(wordforms)**

**Hindi ( hin-x-bible-easy.txt )**
**(before, binding: 5.28 )**

**Hindi ( hin-x-bible-easy.txt )**
**(after, binding: 30.75 )**

Legend (left): ○ यरूशलेम

Legend (middle): ○ जब △ और + वे × फिर

Legend (right): ○ में △ से + के × आये

Levels drawn at 41-46-51%

Levels drawn at 25%

Levels drawn at 25%

26

# Low hanging fruit ...

- Semantic structure of locative case

- Complexity of locative case

- Word order patterns

- **Phylogenetic correlation between word order and complexity**

- Reconstruction of language families based on functional differences

# Density of entropy (only OV languages)



N = 477   Bandwidth = 0.1973

# Phylogenetic modeling



High entropy (1.5+)

Mid entropy (0.1 - 1.5)

Low entropy (0 - 0.1)

VO          OV

# Phylogenetic modeling



High entropy (1.5+)

Mid entropy (0.1 - 1.5)

Low entropy (0 - 0.1)

VO        OV

# Phylogenetic modeling



Average of reconstruction per family

High entropy (1.5+)

Mid entropy (0.1 - 1.5)

Low entropy (0 - 0.1)

VO          OV

# Phylogenetic modeling



Actual frequencies of languages

High entropy (1.5+)

Mid entropy (0.1 - 1.5)

Low entropy (0 - 0.1)

VO          OV

# Phylogenetic modeling



Steady state

High entropy (1.5+)

Mid entropy (0.1 - 1.5)

Low entropy (0 - 0.1)

VO          OV

# Low hanging fruit …

- Semantic structure of locative case

- Complexity of locative case

- Word order patterns

- Phylogenetic correlation between word order and complexity

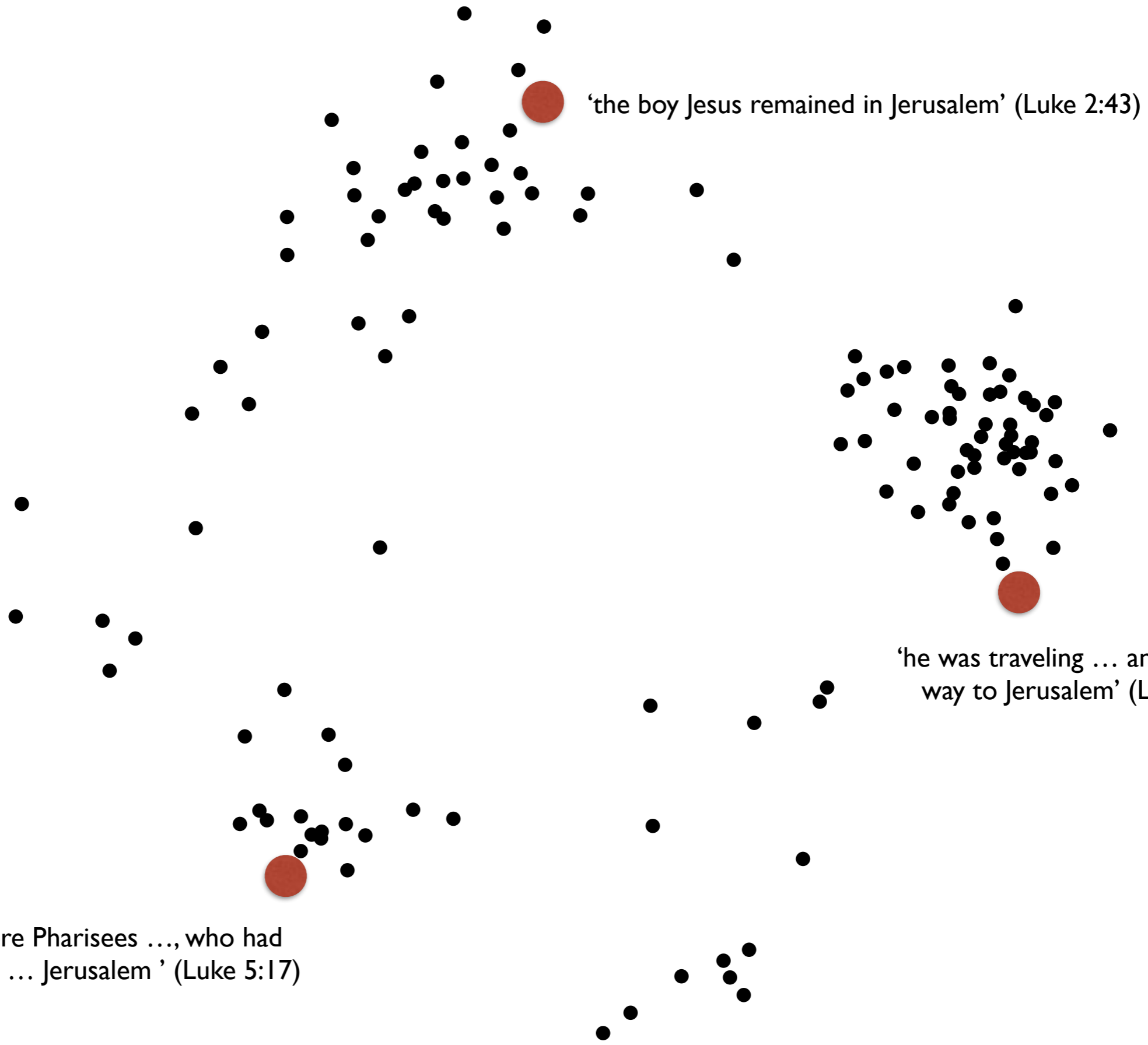- **Reconstruction of language families based on functional differences**

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

- ## Phylogenetic reconstruction

  ‣ Using distribution of cognate affixes over contexts in parallel text (Cysouw & Forker 2009)

  ‣ Simple parsimony ratchet from R-library 'phangorn'

  ‣ Turkic family as an example

Cysouw, Michael & Diana Forker. 2009. Reconstruction of morphosyntactic function: Nonspatial usage of spatial case marking in Tsezic. Language 85(3). 588-617.

chv-x-bible.txt
kjh-x-bible.txt
alt-x-bible.txt
alt-x-bible-southern.txt
gag-x-bible.txt
gag-x-bible-latin.txt
uig-x-bible-romanized.txt
azb-x-bible-latin.txt
tur-x-bible-goodnews.txt
tur-x-bible-southernazeri.txt
tur-x-bible-2009.txt
crh-x-bible.txt
kum-x-bible.txt
krc-x-bible.txt
tat-x-bible.txt
bak-x-bible.txt
kaa-x-bible.txt
kaa-x-bible-latin.txt
kaa-x-bible-cyrillic.txt
tur-x-bible-newworld.txt
kaz-x-bible.txt
kaz-x-bible-newworld1984.txt
kir-x-bible-newworld.txt
kir-x-bible-rayofhope.txt
kir-x-bible-2005.txt

ge
de
den
ø
in
alla
other

Levels drawn at 50%

Levels drawn at 50%

Legend:
- ge (red circle)
- de (green triangle)
- den (orange plus)
- ø (blue x)
- di (magenta diamond)
- alla (yellow inverted triangle)
- other (grey square)

Tree labels:
- chv–x–bible.txt
- kjh–x–bible.txt
- alt–x–bible.txt
- alt–x–bible–southern.txt
- gag–x–bible.txt
- gag–x–bible–latin.txt
- uig–x–bible–romanized.txt
- azb–x–bible–latin.txt
- tur–x–bible–goodnews.txt
- tur–x–bible–southernazeri.txt
- tur–x–bible–2009.txt
- crh–x–bible.txt
- kum–x–bible.txt
- krc–x–bible.txt
- tat–x–bible.txt
- bak–x–bible.txt
- kaa–x–bible.txt
- kaa–x–bible–latin.txt
- kaa–x–bible–cyrillic.txt
- tur–x–bible–newworld.txt
- kaz–x–bible.txt
- kaz–x–bible–newworld1984.txt
- kir–x–bible–newworld.txt
- kir–x–bible–rayofhope.txt
- kir–x–bible–2005.txt

chv−x−bible.txt
kjh−x−bible.txt
alt−x−bible.txt
alt−x−bible−southern.txt
gag−x−bible.txt
gag−x−bible−latin.txt
uig−x−bible−romanized.txt
azb−x−bible−latin.txt
tur−x−bible−goodnews.txt
tur−x−bible−southernazeri.txt
tur−x−bible−2009.txt
crh−x−bible.txt
kum−x−bible.txt
krc−x−bible.txt
tat−x−bible.txt
bak−x−bible.txt
kaa−x−bible.txt
kaa−x−bible−latin.txt
kaa−x−bible−cyrillic.txt
tur−x−bible−newworld.txt
kaz−x−bible.txt
kaz−x−bible−newworld1984.txt
kir−x−bible−newworld.txt
kir−x−bible−rayofhope.txt
kir−x−bible−2005.txt

Legend:
- ○ ge
- △ de
- + den
- ✕ ø
- ◇ in
- ▽ alla
- □ other

Levels drawn at 50%

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

Levels drawn at 50%

ge
ø
de
den
degi
in

Levels drawn at 50%

chv−x−bible.txt
kjh−x−bible.txt
alt−x−bible.txt
alt−x−bible−southern.txt
gag−x−bible.txt
gag−x−bible−latin.txt
uig−x−bible−romanized.txt
azb−x−bible−latin.txt
tur−x−bible−goodnews.txt
tur−x−bible−southernazeri.txt
tur−x−bible−2009.txt
crh−x−bible.txt
kum−x−bible.txt
krc−x−bible.txt
tat−x−bible.txt
bak−x−bible.txt
kaa−x−bible.txt
kaa−x−bible−latin.txt
kaa−x−bible−cyrillic.txt
tur−x−bible−newworld.txt
kaz−x−bible.txt
kaz−x−bible−newworld1984.txt
kir−x−bible−newworld.txt
kir−x−bible−rayofhope.txt
kir−x−bible−2005.txt

ø
ge
de
den
degi
ning

chv−x−bible.txt
kjh−x−bible.txt
alt−x−bible.txt
alt−x−bible−southern.txt
gag−x−bible.txt
gag−x−bible−latin.txt
uig−x−bible−romanized.txt
azb−x−bible−latin.txt
tur−x−bible−goodnews.txt
tur−x−bible−southernazeri.txt
tur−x−bible−2009.txt
crh−x−bible.txt
kum−x−bible.txt
krc−x−bible.txt
tat−x−bible.txt
bak−x−bible.txt
kaa−x−bible.txt
kaa−x−bible−latin.txt
kaa−x−bible−cyrillic.txt
tur−x−bible−newworld.txt
kaz−x−bible.txt
kaz−x−bible−newworld1984.txt
kir−x−bible−newworld.txt
kir−x−bible−rayofhope.txt
kir−x−bible−2005.txt

Levels drawn at 50%

zer
de
den
degi
ning
ø
other

Levels drawn at 50%

Legend:
- ø (blue circle)
- ge (red triangle)
- de (green plus)
- den (yellow x)
- degi (purple diamond)
- ning (orange inverted triangle)

Tree labels:
- chv–x–bible.txt
- kjh–x–bible.txt
- alt–x–bible.txt
- alt–x–bible–southern.txt
- gag–x–bible.txt
- gag–x–bible–latin.txt
- uig–x–bible–romanized.txt
- azb–x–bible–latin.txt
- tur–x–bible–goodnews.txt
- tur–x–bible–southernazeri.txt
- tur–x–bible–2009.txt
- crh–x–bible.txt
- kum–x–bible.txt
- krc–x–bible.txt
- tat–x–bible.txt
- bak–x–bible.txt
- kaa–x–bible.txt
- kaa–x–bible–latin.txt
- kaa–x–bible–cyrillic.txt
- tur–x–bible–newworld.txt
- kaz–x–bible.txt
- kaz–x–bible–newworld1984.txt
- kir–x–bible–newworld.txt
- kir–x–bible–rayofhope.txt
- kir–x–bible–2005.txt

chv-x-bible.txt
kjh-x-bible.txt
alt-x-bible.txt
alt-x-bible-southern.txt
gag-x-bible.txt
gag-x-bible-latin.txt
uig-x-bible-romanized.txt
azb-x-bible-latin.txt
tur-x-bible-goodnews.txt
tur-x-bible-southernazeri.txt
tur-x-bible-2009.txt
crh-x-bible.txt
kum-x-bible.txt
krc-x-bible.txt
tat-x-bible.txt
bak-x-bible.txt
kaa-x-bible.txt
kaa-x-bible-latin.txt
kaa-x-bible-cyrillic.txt
tur-x-bible-newworld.txt
kaz-x-bible.txt
kaz-x-bible-newworld1984.txt
kir-x-bible-newworld.txt
kir-x-bible-rayofhope.txt
kir-x-bible-2005.txt

○ ge
△ de
+ ø
✕ den
◇ di
▽ ning
□ other

Levels drawn at 50%

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

| | |
|---|---|
| ○ | ge |
| △ | de |
| + | ø |
| ✕ | den |
| ◇ | di |
| ▽ | ning |
| □ | other |

Levels drawn at 50%

Levels drawn at 50%

Legend:
- ⭘ ∅
- △ ge
- ✚ de
- ✕ den
- ◇ degi
- ▽ ning

Tree labels:
- chv–x–bible.txt
- kjh–x–bible.txt
- alt–x–bible.txt
- alt–x–bible–southern.txt
- gag–x–bible.txt
- gag–x–bible–latin.txt
- uig–x–bible–romanized.txt
- azb–x–bible–latin.txt
- tur–x–bible–goodnews.txt
- tur–x–bible–southernazeri.txt
- tur–x–bible–2009.txt
- crh–x–bible.txt
- kum–x–bible.txt
- krc–x–bible.txt
- tat–x–bible.txt
- bak–x–bible.txt
- kaa–x–bible.txt
- kaa–x–bible–latin.txt
- kaa–x–bible–cyrillic.txt
- tur–x–bible–newworld.txt
- kaz–x–bible.txt
- kaz–x–bible–newworld1984.txt
- kir–x–bible–newworld.txt
- kir–x–bible–rayofhope.txt
- kir–x–bible–2005.txt

Levels drawn at 50%

chv-x-bible.txt
kjh-x-bible.txt
alt-x-bible.txt
alt-x-bible-southern.txt
gag-x-bible.txt
gag-x-bible-latin.txt
uig-x-bible-romanized.txt
azb-x-bible-latin.txt
tur-x-bible-goodnews.txt
tur-x-bible-southernazeri.txt
tur-x-bible-2009.txt
crh-x-bible.txt
kum-x-bible.txt
krc-x-bible.txt
tat-x-bible.txt
bak-x-bible.txt
kaa-x-bible.txt
kaa-x-bible-latin.txt
kaa-x-bible-cyrillic.txt
tur-x-bible-newworld.txt
kaz-x-bible.txt
kaz-x-bible-newworld1984.txt
kir-x-bible-newworld.txt
kir-x-bible-rayofhope.txt
kir-x-bible-2005.txt

ge
ø
de
den
degi
in

Levels drawn at 50%

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

| | |
|---|---|
| ○ | ge |
| △ | de |
| + | ø |
| ✕ | den |
| ◇ | degi |
| ▽ | i |
| □ | other |

Levels drawn at 50%

Legend:
- ge (red circle)
- de (green triangle)
- ø (blue plus)
- den (yellow x)
- degi (purple diamond)
- i (purple inverted triangle)
- other (gray square)

Tree labels:
- chv-x-bible.txt
- kjh-x-bible.txt
- alt-x-bible.txt
- alt-x-bible-southern.txt
- gag-x-bible.txt
- gag-x-bible-latin.txt
- uig-x-bible-romanized.txt
- azb-x-bible-latin.txt
- tur-x-bible-goodnews.txt
- tur-x-bible-southernazeri.txt
- tur-x-bible-2009.txt
- crh-x-bible.txt
- kum-x-bible.txt
- krc-x-bible.txt
- tat-x-bible.txt
- bak-x-bible.txt
- kaa-x-bible.txt
- kaa-x-bible-latin.txt
- kaa-x-bible-cyrillic.txt
- tur-x-bible-newworld.txt
- kaz-x-bible.txt
- kaz-x-bible-newworld1984.txt
- kir-x-bible-newworld.txt
- kir-x-bible-rayofhope.txt
- kir-x-bible-2005.txt

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

| | |
|---|---|
| ○ | ge |
| △ | de |
| + | ø |
| ✕ | den |
| ◇ | degi |
| ▽ | i |
| ▢ | other |

Levels drawn at 50%

Levels drawn at 50%

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

| | |
|---|---|
| 🟥 | ge |
| 🟩 | de |
| 🟦 | ø |
| 🟨 | den |
| 🟪 | degi |
| 🟦 | ler |
| ⬜ | other |

chv−x−bible.txt
kjh−x−bible.txt
alt−x−bible.txt
alt−x−bible−southern.txt
gag−x−bible.txt
gag−x−bible−latin.txt
uig−x−bible−romanized.txt
azb−x−bible−latin.txt
tur−x−bible−goodnews.txt
tur−x−bible−southernazeri.txt
tur−x−bible−2009.txt
crh−x−bible.txt
kum−x−bible.txt
krc−x−bible.txt
tat−x−bible.txt
bak−x−bible.txt
kaa−x−bible.txt
kaa−x−bible−latin.txt
kaa−x−bible−cyrillic.txt
tur−x−bible−newworld.txt
kaz−x−bible.txt
kaz−x−bible−newworld1984.txt
kir−x−bible−newworld.txt
kir−x−bible−rayofhope.txt
kir−x−bible−2005.txt

ge
de
ø
den
degi
di
other

Levels drawn at 50%

chv-x-bible.txt
kjh-x-bible.txt
alt-x-bible.txt
alt-x-bible-southern.txt
gag-x-bible.txt
gag-x-bible-latin.txt
uig-x-bible-romanized.txt
azb-x-bible-latin.txt
tur-x-bible-goodnews.txt
tur-x-bible-southernazeri.txt
tur-x-bible-2009.txt
crh-x-bible.txt
kum-x-bible.txt
krc-x-bible.txt
tat-x-bible.txt
bak-x-bible.txt
kaa-x-bible.txt
kaa-x-bible-latin.txt
kaa-x-bible-cyrillic.txt
tur-x-bible-newworld.txt
kaz-x-bible.txt
kaz-x-bible-newworld1984.txt
kir-x-bible-newworld.txt
kir-x-bible-rayofhope.txt
kir-x-bible-2005.txt

Legend:
- ge (red circle)
- de (green triangle)
- ø (blue plus)
- den (yellow cross)
- degi (purple diamond)
- di (magenta inverted triangle)

Levels drawn at 50%

chv−x−bible.txt
kjh−x−bible.txt
alt−x−bible.txt
alt−x−bible−southern.txt
gag−x−bible.txt
gag−x−bible−latin.txt
uig−x−bible−romanized.txt
azb−x−bible−latin.txt
tur−x−bible−goodnews.txt
tur−x−bible−southernazeri.txt
tur−x−bible−2009.txt
crh−x−bible.txt
kum−x−bible.txt
krc−x−bible.txt
tat−x−bible.txt
bak−x−bible.txt
kaa−x−bible.txt
kaa−x−bible−latin.txt
kaa−x−bible−cyrillic.txt
tur−x−bible−newworld.txt
kaz−x−bible.txt
kaz−x−bible−newworld1984.txt
kir−x−bible−newworld.txt
kir−x−bible−rayofhope.txt
kir−x−bible−2005.txt

ge
de
ø
den
degi
di

Levels drawn at 50%

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

Levels drawn at 50%

| | |
|---|---|
| ge | (red) |
| de | (green) |
| den | (yellow) |
| ø | (blue) |
| degi | (purple) |
| di | (magenta) |
| other | (gray) |

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

Levels drawn at 50%

ge
de
den
ø
degi
di
other

chv–x–bible.txt
kjh–x–bible.txt
alt–x–bible.txt
alt–x–bible–southern.txt
gag–x–bible.txt
gag–x–bible–latin.txt
uig–x–bible–romanized.txt
azb–x–bible–latin.txt
tur–x–bible–goodnews.txt
tur–x–bible–southernazeri.txt
tur–x–bible–2009.txt
crh–x–bible.txt
kum–x–bible.txt
krc–x–bible.txt
tat–x–bible.txt
bak–x–bible.txt
kaa–x–bible.txt
kaa–x–bible–latin.txt
kaa–x–bible–cyrillic.txt
tur–x–bible–newworld.txt
kaz–x–bible.txt
kaz–x–bible–newworld1984.txt
kir–x–bible–newworld.txt
kir–x–bible–rayofhope.txt
kir–x–bible–2005.txt

ge
de
den
ø
degi
di
other

Levels drawn at 50%

# Outgroup

got–x–bible.txt

bar–x–bible.txt
nld–X–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nob–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt
afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
ltz–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–pattloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

61

got–x–bible.txt
bar–x–bible.txt
nld–x–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nob–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
fry–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–patiloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
deu–x–bible.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

English

62

got–x–bible.txt
bar–x–bible.txt
nld–X–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nob–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
ksh–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–pattloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

German/Dutch/Afrikaans

Danish/Swedish/Norwegian

got–x–bible.txt
bar–x–bible.txt
nld–X–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nob–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt

eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
fry–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–patloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

64

got–x–bible.txt

bar–x–bible.txt
nld–x–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt

???

dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nob–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
fry–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–patiloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

**101**

got-x-bible.txt
bar-x-bible.txt
nld-x-bible-statenvertaling.txt
nld-x-bible-1951.txt
deu-x-bible-erben.txt
deu-x-bible-luther1912.txt
deu-x-bible-luther1545letztehand.txt
deu-x-bible-bolsinger.txt
dan-x-bible-hverdagsdansk.txt
swe-x-bible-nyalevande.txt
nob-x-bible-1988.txt
nno-x-bible-1978.txt
nob-x-bible-1985.txt
nno-x-bible-2011.txt
nno-x-bible-2007.txt
nob-x-bible-1930.txt
nno-x-bible-student.txt
fao-x-bible.txt
isl-x-bible.txt
dan-x-bible-frederik.txt
dan-x-bible-1931.txt
swe-x-bible-newworld.txt
nob-x-bible-newworld.txt
dan-x-bible-newworld.txt
eng-x-bible-contemporary.txt
eng-x-bible-godsword.txt
eng-x-bible-worldwide.txt
eng-x-bible-goodnews.txt
eng-x-bible-newliving.txt
eng-x-bible-newcentury.txt
eng-x-bible-easytoread.txt
eng-x-bible-common.txt
eng-x-bible-newreaders.txt
eng-x-bible-newinternational.txt
eng-x-bible-riverside.txt
eng-x-bible-lexham.txt
eng-x-bible-treeoflife.txt
eng-x-bible-new2007.txt
eng-x-bible-diaglot.txt
eng-x-bible-literal.txt
eng-x-bible-darby.txt
eng-x-bible-newworld2013.txt
eng-x-bible-newworld1984.txt
eng-x-bible-majority.txt
eng-x-bible-world.txt
eng-x-bible-standard.txt
eng-x-bible-scriptures.txt
eng-x-bible-newsimplified.txt
eng-x-bible-etheridge.txt
eng-x-bible-basic.txt
enm-x-bible-wycliffe.txt
eng-x-bible-catholic.txt
eng-x-bible-montgomery.txt
eng-x-bible-plontz.txt
eng-x-bible-amplified.txt
eng-x-bible-coverdale.txt
enm-x-bible-kingjames.txt
eng-x-bible-geneva.txt
enm-x-bible-tyndale.txt
enm-x-bible-bishop.txt
afr-x-bible-boodskap.txt
afr-x-bible-lewende.txt
nld-x-bible-2007.txt
deu-x-bible-volxbibel.txt
swg-x-bible.txt
nld-x-bible-newworld.txt
nld-x-bible-2004.txt
nld-x-bible.txt
deu-x-bible-textbibel.txt
deu-x-bible-pattloch.txt
deu-x-bible-newworld.txt
deu-x-bible-meister.txt
deu-x-bible-konkordant.txt
deu-x-bible-interlinear.txt
deu-x-bible-albrecht.txt
nds-x-bible.txt
deu-x-bible-freebible.txt
deu-x-bible-elberfelder1905.txt
deu-x-bible-elberfelder1871.txt
afr-x-bible-newworld.txt
afr-x-bible-1953.txt
afr-x-bible-viralmal.txt
afr-x-bible-1983.txt
deu-x-bible-zuercher.txt
pdt-x-bible.txt
deu-x-bible-tafelbibel.txt
deu-x-bible-schlachter2000.txt
deu-x-bible-schlachter.txt
deu-x-bible-gruenewalder.txt
deu-x-bible-menge.txt
deu-x-bible-greber.txt
deu-x-bible-neue.txt
deu-x-bible-hoffnung.txt
deu-x-bible-genfer2011.txt

Legend:
○ in
△ nach
+ aus

66

99

got–x–bible.txt
bar–x–bible.txt
nld–x–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nno–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
eng–x–bible–kingjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
nld–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–pattloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

○ in
△ nach
+ aus
✕ von

**101**

got-x-bible.txt
bar-x-bible.txt
nld-x-bible-statenvertaling.txt
nld-x-bible-1951.txt
deu-x-bible-erben.txt
deu-x-bible-luther1912.txt
deu-x-bible-luther1545letztehand.txt
deu-x-bible-bolsinger.txt
dan-x-bible-hverdagsdansk.txt
swe-x-bible-nyalevande.txt
nob-x-bible-1988.txt
nno-x-bible-1978.txt
nob-x-bible-1985.txt
nno-x-bible-2011.txt
nno-x-bible-2007.txt
nob-x-bible-1930.txt
nno-x-bible-student.txt
fao-x-bible.txt
isl-x-bible.txt
dan-x-bible-frederik.txt
dan-x-bible-1931.txt
swe-x-bible-newworld.txt
nob-x-bible-newworld.txt
dan-x-bible-newworld.txt
eng-x-bible-contemporary.txt
eng-x-bible-godsword.txt
eng-x-bible-worldwide.txt
eng-x-bible-goodnews.txt
eng-x-bible-newliving.txt
eng-x-bible-newcentury.txt
eng-x-bible-easytoread.txt
eng-x-bible-common.txt
eng-x-bible-newreaders.txt
eng-x-bible-newinternational.txt
eng-x-bible-riverside.txt
eng-x-bible-lexham.txt
eng-x-bible-treeoflife.txt
eng-x-bible-new2007.txt
eng-x-bible-diaglot.txt
eng-x-bible-literal.txt
eng-x-bible-darby.txt
eng-x-bible-newworld2013.txt
eng-x-bible-newworld1984.txt
eng-x-bible-majority.txt
eng-x-bible-world.txt
eng-x-bible-standard.txt
eng-x-bible-scriptures.txt
eng-x-bible-newsimplified.txt
eng-x-bible-etheridge.txt
eng-x-bible-basic.txt
enm-x-bible-wycliffe.txt
eng-x-bible-catholic.txt
eng-x-bible-montgomery.txt
eng-x-bible-plontz.txt
eng-x-bible-amplified.txt
enm-x-bible-coverdale.txt
eng-x-bible-kingjames.txt
enm-x-bible-geneva.txt
enm-x-bible-tyndale.txt
enm-x-bible-bishop.txt

afr-x-bible-boodskap.txt
afr-x-bible-lewende.txt
nld-x-bible-2007.txt
deu-x-bible-volxbibel.txt
swg-x-bible.txt
nld-x-bible-newworld.txt
nld-x-bible-2004.txt
nld-x-bible.txt
deu-x-bible-textbibel.txt
deu-x-bible-pattloch.txt
deu-x-bible-newworld.txt
deu-x-bible-meister.txt
deu-x-bible-konkordant.txt
deu-x-bible-interlinear.txt
deu-x-bible-albrecht.txt
nds-x-bible.txt
deu-x-bible-freebible.txt
deu-x-bible-elberfelder1905.txt
deu-x-bible-elberfelder1871.txt
afr-x-bible-newworld.txt
afr-x-bible-1953.txt
afr-x-bible-viralmal.txt
afr-x-bible-1983.txt
deu-x-bible-zuercher.txt
pdt-x-bible.txt
deu-x-bible-tafelbibel.txt
deu-x-bible-schlachter2000.txt
deu-x-bible-schlachter.txt
deu-x-bible-gruenewalder.txt
deu-x-bible-menge.txt
deu-x-bible-greber.txt
deu-x-bible-neue.txt
deu-x-bible-hoffnung.txt
deu-x-bible-genfer2011.txt

Legend:
○ in
△ nach
+ aus

128

got-x-bible.txt
bar-x-bible.txt
nld-x-bible-statenvertaling.txt
nld-x-bible-1951.txt
deu-x-bible-erben.txt
deu-x-bible-luther1912.txt
deu-x-bible-luther1545letztehand.txt
deu-x-bible-bolsinger.txt
dan-x-bible-hverdagsdansk.txt
swe-x-bible-nyalevande.txt
nob-x-bible-1988.txt
nno-x-bible-1978.txt
nob-x-bible-1985.txt
nob-x-bible-2011.txt
nno-x-bible-2011.txt
nno-x-bible-2007.txt
nob-x-bible-1930.txt
nno-x-bible-student.txt
isl-x-bible.txt
fao-x-bible.txt
dan-x-bible-frederik.txt
dan-x-bible-1931.txt
swe-x-bible-newworld.txt
nob-x-bible-newworld.txt
dan-x-bible-newworld.txt

eng-x-bible-contemporary.txt
eng-x-bible-godsword.txt
eng-x-bible-worldwide.txt
eng-x-bible-goodnews.txt
eng-x-bible-newliving.txt
eng-x-bible-newcentury.txt
eng-x-bible-easytoread.txt
eng-x-bible-common.txt
eng-x-bible-newreaders.txt
eng-x-bible-newinternational.txt
eng-x-bible-riverside.txt
eng-x-bible-lexham.txt
eng-x-bible-treeoflife.txt
eng-x-bible-new2007.txt
eng-x-bible-diaglot.txt
eng-x-bible-literal.txt
eng-x-bible-darby.txt
eng-x-bible-newworld2013.txt
eng-x-bible-newworld1984.txt
eng-x-bible-majority.txt
eng-x-bible-world.txt
eng-x-bible-standard.txt
eng-x-bible-scriptures.txt
eng-x-bible-newsimplified.txt
eng-x-bible-etheridge.txt
eng-x-bible-basic.txt
eng-x-bible-wycliffe.txt
eng-x-bible-catholic.txt
eng-x-bible-montgomery.txt
eng-x-bible-clontz.txt
eng-x-bible-amplified.txt
enm-x-bible-coverdale.txt
eng-x-bible-kindjames.txt
enm-x-bible-geneva.txt
enm-x-bible-tyndale.txt
enm-x-bible-bishop.txt

afr-x-bible-boodskap.txt
afr-x-bible-lewende.txt
nld-x-bible-2007.txt
deu-x-bible-volxbibel.txt
swg-x-bible.txt
nld-x-bible-newworld.txt
nld-x-bible-2004.txt
nld-x-bible-x-bible.txt
deu-x-bible-textbibel.txt
deu-x-bible-pattloch.txt
deu-x-bible-newworld.txt
deu-x-bible-meister.txt
deu-x-bible-konkordant.txt
deu-x-bible-interlinear.txt
deu-x-bible-albrecht.txt
nds-x-bible.txt
deu-x-bible-freebible.txt
deu-x-bible-elberfelder1905.txt
deu-x-bible-elberfelder1871.txt
afr-x-bible-newworld.txt
afr-x-bible-1953.txt
afr-x-bible-viralmal.txt
afr-x-bible-1983.txt
deu-x-bible-zuercher.txt
pdt-x-bible.txt
deu-x-bible-tafelbibel.txt
deu-x-bible-schlachter2000.txt
deu-x-bible-schlachter.txt
deu-x-bible-gruenewalder.txt
deu-x-bible-menge.txt
deu-x-bible-greber.txt
deu-x-bible-neue.txt
deu-x-bible-hoffnung.txt
deu-x-bible-genfer2011.txt

Legend:
○ in
△ from
+ back

69

103

got-x-bible.txt
bar-x-bible.txt
nld-x-bible-statenvertaling.txt
nld-x-bible-1951.txt
deu-x-bible-erben.txt
deu-x-bible-luther1912.txt
deu-x-bible-luther1545letztehand.txt
deu-x-bible-boisinger.txt
dan-x-bible-hverdagsdansk.txt
swe-x-bible-nyalevande.txt
nob-x-bible-1988.txt
nno-x-bible-1978.txt
nob-x-bible-1985.txt
nob-x-bible-2011.txt
nno-x-bible-2011.txt
nno-x-bible-2007.txt
nob-x-bible-1930.txt
nno-x-bible-student.txt
isl-x-bible.txt
fao-x-bible.txt
dan-x-bible-frederik.txt
dan-x-bible-1931.txt
swe-x-bible-newworld.txt
nob-x-bible-newworld.txt
dan-x-bible-newworld.txt
eng-x-bible-contemporary.txt
eng-x-bible-godsword.txt
eng-x-bible-worldwide.txt
eng-x-bible-goodnews.txt
eng-x-bible-newliving.txt
eng-x-bible-newcentury.txt
eng-x-bible-easytoread.txt
eng-x-bible-common.txt
eng-x-bible-newreaders.txt
eng-x-bible-newinternational.txt
eng-x-bible-riverside.txt
eng-x-bible-lexham.txt
eng-x-bible-treeoflife.txt
eng-x-bible-new2007.txt
eng-x-bible-diaglot.txt
eng-x-bible-literal.txt
eng-x-bible-darby.txt
eng-x-bible-newworld2013.txt
eng-x-bible-majority.txt
eng-x-bible-newworld1984.txt
eng-x-bible-world.txt
eng-x-bible-standard.txt
eng-x-bible-scriptures.txt
eng-x-bible-newsimplified.txt
eng-x-bible-etheridge.txt
eng-x-bible-basic.txt
enm-x-bible-wycliffe.txt
eng-x-bible-catholic.txt
eng-x-bible-montgomery.txt
eng-x-bible-plontz.txt
eng-x-bible-amplified.txt
eng-x-bible-coverdale.txt
enm-x-bible-kingjames.txt
eng-x-bible-geneva.txt
enm-x-bible-tyndale.txt
enm-x-bible-bishop.txt

afr-x-bible-boodskap.txt
afr-x-bible-lewende.txt
nld-x-bible-2007.txt
deu-x-bible-volxbibel.txt
swg-x-bible.txt
nld-x-bible-newworld.txt
nld-x-bible-2004.txt
nld-x-x-bible.txt
deu-x-bible-textbibel.txt
deu-x-bible-pattloch.txt
deu-x-bible-newworld.txt
deu-x-bible-meister.txt
deu-x-bible-konkordant.txt
deu-x-bible-interlinear.txt
deu-x-bible-albrecht.txt
nds-x-bible.txt
deu-x-bible-freebible.txt
deu-x-bible-elberfelder1905.txt
deu-x-bible-elberfelder1871.txt
afr-x-bible-newworld.txt
afr-x-bible-1953.txt
afr-x-bible-viralmal.txt
afr-x-bible-1983.txt
deu-x-bible-zuercher.txt
pdt-x-bible.txt
deu-x-bible-tafelbibel.txt
deu-x-bible-schlachter2000.txt
deu-x-bible-schlachter.txt
deu-x-bible-gruenewalder.txt
deu-x-bible-menge.txt
deu-x-bible-greber.txt
deu-x-bible-neue.txt
deu-x-bible-hoffnung.txt
deu-x-bible-genfer2011.txt

to
in
from
back
near
ward

70

**137**



got–x–bible.txt
bar–x–bible.txt
nld–x–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–bolsinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nob–x–bible–2011.txt
nno–x–bible–2011.txt
nno–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–treeoflife.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
eng–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–clontz.txt
eng–x–bible–amplified.txt
enm–x–bible–coverdale.txt
enm–x–bible–geneva.txt
enm–x–bible–kindjames.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt
afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
fry–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–pattloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

| | |
|---|---|
| ○ | to |
| △ | in |
| + | aus |

**138**



got–x–bible.txt
bar–x–bible.txt
nld–x–bible–statenvertaling.txt
nld–x–bible–1951.txt
deu–x–bible–erben.txt
deu–x–bible–luther1912.txt
deu–x–bible–luther1545letztehand.txt
deu–x–bible–eisinger.txt
dan–x–bible–hverdagsdansk.txt
swe–x–bible–nyalevande.txt
nob–x–bible–1988.txt
nno–x–bible–1978.txt
nob–x–bible–1985.txt
nno–x–bible–2011.txt
nno–x–bible–2007.txt
nob–x–bible–1930.txt
nno–x–bible–student.txt
isl–x–bible.txt
fao–x–bible.txt
dan–x–bible–frederik.txt
dan–x–bible–1931.txt
swe–x–bible–newworld.txt
nob–x–bible–newworld.txt
dan–x–bible–newworld.txt
eng–x–bible–contemporary.txt
eng–x–bible–godsword.txt
eng–x–bible–worldwide.txt
eng–x–bible–goodnews.txt
eng–x–bible–newliving.txt
eng–x–bible–newcentury.txt
eng–x–bible–easytoread.txt
eng–x–bible–common.txt
eng–x–bible–newreaders.txt
eng–x–bible–newinternational.txt
eng–x–bible–riverside.txt
eng–x–bible–lexham.txt
eng–x–bible–freeonline.txt
eng–x–bible–new2007.txt
eng–x–bible–diaglot.txt
eng–x–bible–literal.txt
eng–x–bible–darby.txt
eng–x–bible–newworld2013.txt
eng–x–bible–newworld1984.txt
eng–x–bible–majority.txt
eng–x–bible–world.txt
eng–x–bible–standard.txt
eng–x–bible–scriptures.txt
eng–x–bible–newsimplified.txt
eng–x–bible–etheridge.txt
eng–x–bible–basic.txt
enm–x–bible–wycliffe.txt
nld–x–bible–catholic.txt
eng–x–bible–montgomery.txt
eng–x–bible–plontz.txt
eng–x–bible–amplified.txt
eng–x–bible–coverdale.txt
enm–x–bible–kindjames.txt
enm–x–bible–geneva.txt
enm–x–bible–tyndale.txt
enm–x–bible–bishop.txt

afr–x–bible–boodskap.txt
afr–x–bible–lewende.txt
nld–x–bible–2007.txt
deu–x–bible–volxbibel.txt
swg–x–bible.txt
nld–x–bible–newworld.txt
nld–x–bible–2004.txt
nld–x–bible.txt
deu–x–bible–textbibel.txt
deu–x–bible–pattloch.txt
deu–x–bible–newworld.txt
deu–x–bible–meister.txt
deu–x–bible–konkordant.txt
deu–x–bible–interlinear.txt
deu–x–bible–albrecht.txt
nds–x–bible.txt
deu–x–bible–freebible.txt
deu–x–bible–elberfelder1905.txt
deu–x–bible–elberfelder1871.txt
afr–x–bible–newworld.txt
afr–x–bible–1953.txt
afr–x–bible–viralmal.txt
afr–x–bible–1983.txt
deu–x–bible–zuercher.txt
pdt–x–bible.txt
deu–x–bible–tafelbibel.txt
deu–x–bible–schlachter2000.txt
deu–x–bible–schlachter.txt
deu–x–bible–gruenewalder.txt
deu–x–bible–menge.txt
deu–x–bible–greber.txt
deu–x–bible–neue.txt
deu–x–bible–hoffnung.txt
deu–x–bible–genfer2011.txt

| | |
|---|---|
| ○ | to |
| △ | von |