### Extracting relational data from Wikipedia

by Χρήστος Περιβολαρόπουλος

Μια διπλοματική διατφιβη.

Επιβλέποντες καθηγητές: Κυριάκος Σγάρμπας, Boris Katz

University of Patras

June 2015

#### Χρήστος Περιβολαρόπουλος

#### **Abstract**

Το START (SynTactic Analysis using Reversible Transformations) είναι ένα παμέτο λογισμικού γραμμένο σε γλώσσα Common Lisp το οποίο αντλεί πληροφορίες από διαδικτυακές πηγές και τις χρησιμοποιεί για απαντήσει σε αυθαίρετες ερωτήσεις που δέχεται. Αναπτύχθηκε στο εργαστήριο Infolab του ΜΙΤ. Για τον εμπλουτισμό των πληροφοριών που χρησιμοποιεί το START χρησιμοποιείται το πακέτο λογισμικού Omnibase μέσω του οποίου επιτυγχάνεται πρόσβαση του START σε πολλαπλές πήγες στο διαδίκτυο. Στην παρούσα εργασία αναπτύσσουμε μια επέκταση του Omnibase που επιτρέπει την πρόσβαση του START στην wikipedia και στις πληροφορίες που αυτή περιέχει. Η επέκταση αυτή του Omnibase ονομάζεται wikipediabase. Για την ευκολότερη πρόσβαση του START στη Wikipedia δημιουργήσαμε ένα πρόγραμμα (wikipedia-mirror) που δημιουργεί κλώνους της Wikipedia που αποθηκεύονται τοπικά ανεξάρτητα από το διαδίκτυο. Έτσι επιτυγχάνεται ταχύτερη και πιο αξιόπιστη πρόσβαση του START στο σύνολο των δεδομένων της wikipedia.

START (SynTactic Analysis using Reversible Transformations) is a piece of sofware written in common lisp that retrieves information from internet resources and uses them to answer to arbitrary natural language questions. It was developed in InfoLab of MiT. For the enrichment of the retrieved information it uses the Omnibase software through which START gets access to multiple sources on the internet. In the present thesis we present an extension to Omnibase that allows START to get access to wikipedia and the information that it contains. This extension to Omnibase is called WikipediaBase. For easier access to wikipedia we also developed a separate program (wikipedia-mirror) that creates clones of wikipedia running locally and independently to the internet. This way faster and more reliable access to wikipedia is accomplished.

"The Initial Mystery that attends any journey is how did the traveler reach his starting point in the first place?"

- Louise Bogan, Journey Around My Room

### Acknowledgements

First and foremost, I have to thank my research supervisor Professor L. Perivolaropoulos. Without his assistance and dedicated involvement in every step of the way, this thesis would have never been accomplished.

I also take this opportunity to express my gratitude to all the Department faculty members, for their help and support throughout my years of studying.

Lastly, I offer my sincere thanks to one and all who directly and indirectly have lent a helping hand in this virtue.

# **Contents**

# Part I Εισαγωγή

Το START Natural Language System[1] είναι ένα σύστημα λογισμικού που έχει σχεδιαστεί για να απαντά σε ερωτήσεις που τίθενται σε αυτό σε φυσική γλώσσα. Το START σαρώνει την εισερχόμενη ερώτηση, αναλύει τη δομή της και προσπαθεί να εντοπίσει παρόμοιες δομές στη γνωστική του βάση δεδομένων με στόχο να βρει απάντηση στην ερώτηση που του έχει τεθεί. Με τον τρόπο αυτό, το START παρέχει σε μη εξειδηκευμένους χρήστες γρήγορη πρόσβαση σε γνώση, η ευρεση της οποίας σε πολλές περιπτώσεις μπορεί να είναι χρονοβόρα ακόμα και για τον ειδικό.

Το START δημιουργήθηκε από τον Dr Boris Katz στο Artificial Intelligence Laboratory του MIT. Επί του παρόντος, το σύστημα υφίσταται περαιτέρω ανάπτυξη από το InfoLab Group, με επικεφαλή τον Dr Katz. Το START για πρώτη φορά συνδέθηκε με το World Wide Web το Δεκέμβριο του 1992, και με τις διάφορες μορφές του έχει μέχρι σήμερα απαντήσει σε εκατομμύρια ερωτήσεις χρηστών από όλο τον κόσμο.

Το START μποφεί να χειφιστεί μεγάλη ποιχιλία μέσων, συμπεφιλαμβανομένων χειμένων, διαγφαμμάτων, ειχόνων, βίντεο και ήχων, ιστοσελίδων και άλλων χφησιμοποιώντας μια βασική τεχνική που ονομάζεται "natural language annotation". Κατά την τεχνική αυτή παίφνουμε ως είσοδο μια αντιστοιχία ανάμεσα σε αυθαίφετα δεδομένα και φφάσεις ή πφοτάσεις που τα πεφιγφάφουν και χφησιμοποιούμε τις φφάσεις ή πφοτάσεις αυτές για να απαντηθούν οι εφωτήσεις του χρήστη με τα αντίστοιχα δεδομένα.

Η επεξεργασία της φυσικής γλώσσας του START αποτελείται από δύο μονάδες που αναγνωρίζουν την ίδια σύνταξη. Αυτές οι δυο μονάδες είναι: Η μονάδα κατανόησης και η μονάδα παραγωγής γλώσσας. Η μονάδα κατανόησης αναλύει το αγγλικό κείμενο και παράγει μια γνωστική βάση που κωδικοποιεί τις πληροφορίες του κειμένου. Η μονάδα παραγωγής γλώσσας παράγει αγγλικές προτάσεις λαμβάνοντας υπόψη μόνο ένα κατάλληλο τμήμα της γνωστικής βάσης. Αυτές οι μονάδες σε συνδυασμό με την τεχνική της "natural language annotation", δίνουν την δυνατότητα της παραγωγής απαντήσεων σε φυσική γλώσσα.

Το Omnibase[2] είναι μια "ειχονιχή" βάση δεδομένων που το START χρησιμοποιεί για να έχει ομοιογενή πρόσβαση σε ετερογενείς πήγες γνώσης. Για παράδειγμα δίνει ομογενή πρόσβαση στο imdb[?] και στο CIA world factbook[?]. Το Omnibase αναπτύχθηκε για πρώτη φορά το 2002, περίπου το ίδιο χρονικό διάστημα που η wikipedia έκανε την πρώτη της εμφάνιση (2001).

Η διαδικτυακή εγκυκλοπαίδεια Wikipedia[?] είναι ένα τεφάστιο, συνεχώς εξελισσόμενο δίκτυο από αλληλένδετες πληφοφοφίες σε μοφφή κειμένου. Οι πληφοφοφίες αυτές εξελίσσονται και εμπλουτίζονται χειφοκίνητα από τους ίδιους τους χφήστες. Αποτελεί έναν απαφάμιλλο και σε μεγάλο βαθμό ανεκμετάλλευτο πόφο για την επεξεργασία φυσικής γλώσσας, διαχείφιση

της γνώσης, εξόουξη δεδομένων, και διάφορους άλλους τομείς έρευνας. Είναι το προϊόν της συνεργασίας εκατομμυρίων ανθρώπων. Η Wikipedia βασίζεται στο σύστημα wiki. Το σύστημα αυτό αποτελεί μια κατηγορία ιστοσελίδων που επιτρέπουν την συνεργατική τροποποίηση περιεχομένου τους από τους χρήστες.

Λόγω της πολυπλοκότητας και της ιδιαίτερα αδόμητης φύσης της wikipedia, αντί του omnibase αναπτύξαμε μια ξεχωριστή υπηρεσία, το WikipediaBase, για να την καταστήσουμε προσβάσιμη από το START. Το wikipediaBase αποτελεί και το αντικείμενο της παρούσας διατριβής. Επιπλέον, για να αποφευχθεί ο κορεσμός του wikiedpedia.org, δημιουργήσαμε το wikipedia-mirror, που δημιουργεί κλώνους της wikipedia.

Σε πολλά άρθοα της Wikipedia συναντάμε μια δομή, το infobox, που χρησιμοποιείται συχνά από το WikipediaBase. Τα Infoboxes είναι πίνακες που χρησιμοποιούνται συνήθως στη wikipedia για να παρέχουν την περίληψη ενός άρθου με ήμι δομημένο τρόπο. Τα Infoboxes είναι η κύρια πηγή πληροφοριών για τη WikipediaBase

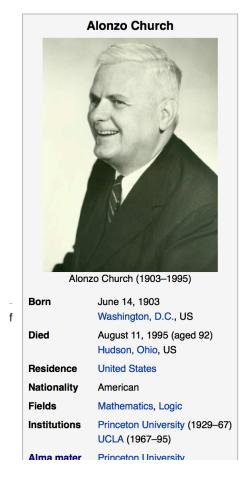


Figure 1: Ένα παράδειγμα ενός infobox

Σε οφους mediawiki markup, ένα infobox είναι ένα typed template που αποδίδεται σε html. Για παφάδειγμα:

```
{{Infobox scientist
| name
                    = Gerhard Gentzen
| image
                    = Gerhard Gentzen.jpg
| image_size
| alt
| caption
                    = Gerhard Gentzen in Prague, 1945.
                  = \{\{Birth date | 1909 | 11 | 24\}\}
| birth_date
| birth_place
                  = [[Greifswald]], [[Germany]]
                   = {{Death date and age
| death_date
   |1945|8|4|1909|11|24}}
```

#### Θα παράξει το εξής:

Οι τύποι του Infobox, αναφερόμενοι και ως κλάσεις, είναι οργανωμένοι με μια αρκετά ευρεία ιεραρχία[9]. Για παράδειγμα Template: Infobox Austrian district είναι μια ειδική περίπτωση ενός Template: Infobox settlement και το καθένα μετατρέπεται από mediawiki markup σε HTML (rendering) με διαφορετικό τρόπο.

Ένα άφθοο μπορεί να έχει περισσότερα από ένα infoboxes, για παράδειγμα, το άρθοο για τον Bill Clinton έχει δύο infobox: ένα για Infobox Officeholder και ένα για Infobox President.

# Part II Wikipediabase

Η WikipediaBase είναι μια πηγή πληφοφοφιών χφησιμοποιούμενη από το START. Είναι υπεύθυνη για την παφοχή πφόσβασης σε πληφοφοφίες που σχετίζονται με την wikipedia. Μιμείται τον τφόπο επικοινωνίας του το Omnibase. Μέτα την αφχική της έκδοση η WikipediaBase έχει ξαναγφαφεί δυο φοφές. Η αφχική έκδοση ήταν γφαμμένη σε Java. Στη συνέχεια ξαναγφάφτηκε σε Ruby διατηφώντας την αφχική αφχιτεκτονική και το σχεδιασμό, και η παφούσα έκδοση ως αντικείμενο της παφούσας διατφιβής είναι γφαμμένη σε σε python με νέο σχεδιασμό και αφχιτεκτονική.

Υπάρχουν δύο βασικοί λόγοι για την τελευταία επανεγγραφή: η Python επελέγη διότι διδάσκεται ως προπτυχιακό μάθημα στο ΜΙΤ, και ως εκ τούτου, ένα πρόγραμμα σε Python θα κάνει την εισαγωγή νέων φοιτητών του ΜΙΤ στην ομάδα εργασίας ομαλότερη. Ο δεύτερος και πιο σημαντικός όμως λόγος για την επανεγγραφή είναι ότι ενώ ο αρχικός σχεδιασμός του προηγούμενου WikipediaBase ήταν στην αρχή επαρκής, στη συνέχεια η WikipediaBase μεγάλωσε σε τέτοιο σημείο, όπου ο κώδικας περιείχε πολλές ειδικές περιπτώσεις και ήταν δύσκολο να κατανοηθεί.

Το WikipediaBase σε python αρχικά γράφτηκε από τον Χρήστο Περιβολαρόπουλο σε στενή συνεργασία με την Dr Sue Felshin και τελικά παραδόθηκε στους Sue Felshin, Alvaro Morales και ton Michael Silver. Αργότερα και άλλοι φοιτητές εντάχθηκαν στο έργο.

### Chapter 1

# Λειτουργικότητα

Στην επικοινωνία με το START υπεισέρχονται κάποιες βασικές έννοιες:

- symbol είναι μια έννοια στην οποία μπορούμε να αναφερθούμε.
- attribute είναι ένα χαρακτηριστικά ενός symbol.
- class είναι σύνολα που περιέχουν διάφορες έννοιες. Κάθε έννοια μπορεί να εντάσσεται σε παραπάνω από ένα class. Όλα τα symbols που εντάσσονται σε ένα class υποστηρίζουν ένα συγκεκριμένο σύνολο από attributes.

Υπάρχει ένα προς ένα αντιστοιχία αυτών των εννοιών με τις βασικές έννοιες του infobox. Δηλαδή το άρθρο στο οποίο αναφέρεται ένα infobox είναι ένα symbol, το class αυτού του symbol είναι η κλάση του infobox και τα χαρακτηριστικά που εκφράζει το infobox είναι τα attributes.

Όλα τα αντιπείμενα της WikipediaBase ανήπουν πληφονομιπά στην υπεφπλάση wikibase-term, η οποία υποστηφίζει τα χαφαπτηφιστιπά IMAGE-DATA, SHORT-ARTICLE, URL, COORDINATES, PROPER, παι NUMBER.

Οι εντολές της WikipediaBase και οι τιμές επιστροφής τους χρησιμοποιούν κωδικοποίηση s-expressions. s-expressions είναι εκφράσεις που έχουν τη μορφή (a b c (d e) f g).

Η WikipediaBase παρέχει τις απόλουθες λειτουργίες:

#### H εντολή get

Δεδομένης μιας class, ενός symbol, και ενός τυποποιημένου attribute, δηλαδή ενός χαρακτηριστικού με typecode, η εντολή get επιστρέφει την τιμή του attribute. Έγκυρα typecodes χαρακτηριστικών είναι : code

(για ένα attribute όνομα όπως στο infobox wiki markup) και :rendered (για ένα attribute όνομα στο rendered form από το infobox).

#### (a) Typecodes

Οι get εντολές πρέπει να επιστρέφουν μια λίστα από τυποποιημένες τιμές, δηλαδή ενα ζευγάρι τιμής - typecode. Έγκυρα typecodes είναι:

#### i. :HTML

Μια συμβολοσειρά προσαρμοσμένη για μετατροπή σε HTML. Η συμβολοσειρά πρέπει να είναι escaped για lisp, εννοώντας quoted, και με double quotes και backslashes escaped με backslashes. Η συμβολοσειρά δεν απαιτείται να περιέχει HTML κώδικες. Για παράδειγμα:

```
(get "wikipedia-sea" "Black_Sea" (:code "
    AREA"))
=> ((:html "436,402_km2_(168,500_sq_mi)"))

(get "wikipedia-president" "Bill_Clinton" (:
    code "SUCCESSOR"))
=> ((:html "George_W._Bush"))

(get "wikipedia-president" "Bill_Clinton" (:
    rendered "Succeeded_by"))
=> ((:html "George_W._Bush"))
```

#### ii. : YYYYMMDD

Οι αναλυμένες ημερομηνίες αντιπροσωπεύονται σαν αριθμοί, χρησιμοποιώντας τον τύπο ΥΥΥΥΜΝΟΟ με αρνητικούς αριθμούς αντιπροσωπεύονται οι  $\Pi.X.$  ημερομηνίες.

(Οι μη αναλυμένες ημεφομηνίες αντιπφοσωπεύονται σαν HTML strings χρησιμοποιώντας το : HTML typecode.)

```
(get "wikibase-person" "Barack Obama" (:ID "
    BIRTH-DATE"))
=> ((:yyyymmdd 19610804))

(get "wikibase-person" "Julius Caesar" (:ID
    "BIRTH-DATE"))
=> ((:YYYYMMDD -1000713))
```

#### iii. : CALCULATED

Το Typecode για χαρακτηριστικά υπολογισμένα με βάση χαρακτηριστικά του άρθρου, πχ., GENDER and NUMBER. Βλέπε παρακάτω στο Special Attributes για την ολοκληρωμένη λίστα των υπολογισμένων attributes.

iv. : CODE

Ξεπερασμένο συνώνυμο του : HTML.

v. : STRING

Ξεπερασμένο συνώνυμο του : HTML.

vi. Special Attributes

Μερικά χαρακτηριστικά είναι ειδικά επειδή υπολογίζονται από τη WikipediaBase αντί να προέρχονται από infoboxes. Αυτά τα χαρακτηριστικά θα πρέπει να είναι ειδικά για τις classes wikibase-term, wikibase-person, και wikipedia-paragraphs.

- A. SHORT-ARTICLE, για την class wikibase-term Η πρώτη παράγραφος του άρθρου. Αν η πρώτη παράγραφος είναι μικρότερη από 350 χαρακτήρες, τότε η επιστρεφόμενη τιμή είναι το πρώτο μέρος του κειμένου έτσι ώστε το άθροισμα των χαρακτήρων είναι τουλάχιστον 350.
- B. URL, για την class wikibase-term Το URL του άρθρου ως ((:url URL))
- C. IMAGE-DATA, για την class wikibase-term
  Μια λίστα από URLs εικόνων στο περιεχόμενο του άρθρου
  (αποκλείει εικόνες που είναι στη σελίδα αλλά εκτός του
  περιεχομένου του άρθρου). Εάν δεν υπάρχουν εικόνες
  επιστρέφει μια κενή λίστα.
  Η "καλύτερη" εικόνα πρέπει να είναι η πρώτη της λίστας,
  εάν υπάρχει εικόνα στην κορυφή του infobox, αυτή θεωρείται
  η καλύτερη εικόνα, διαφορετικά είναι η πρώτη εικόνα
  που εμφανίζεται οπουδήποτε στο άρθρο. Εαν δεν υπάρχει
  caption, η τιμή του caption παραλείπεται
  π.χ.,προτιμότερο ((0 "Harimau\\_Harimau\\_cover.jpg"))
  από ((0 "Harimau\\_Harimau\\_cover.jpg" "")).
- D. COORDINATES, για την class wikibase-term Το γεωγραφικό πλάτος και μήκος. Εντοπίζονται είτε στο πάνω δεξιό άκρο του άρθρου, είτε στο infobox. Η τιμή είναι μια λίστα του πλάτους και μήκους, πχ. ((:coordinates latitude longitude))

Figure 1.1: An example of coordinates in the header

- Ε. ΒΙRTH-DATE, για την class wikibase-person
  Η ημεφομηνία γέννησης. Λαμβάνεται από το infobox, το άφθρο, ή τις πληφοφορίες κατηγορίας του άφθρου.
  Η τιμή μποφεί να είναι μια a parsed or unparsed date. Οι parsed dates αντιπροσωπεύονται ως αφιθμοί, χρησιμοποιώντας τη μοφφή ΥΥΥΥΜΜDD.
- F. DEATH-DATE, για την class wikibase-person
  Η ημεφομηνία θανάτου. Λαμβάνεται με παφόμοιο τφόπο
  όπως το BIRTH-DATE. Επιστφέφει τον ίδιο τύπο τιμής όπως
  ΒΙRTH-DATE, εκτός αν το πφόσωπο ζει, τότε βγάζει διευκφίνηση
  (error /"Currently alive"/).
- G. GENDER, για την class wikibase-person
  Το φύλο του προσώπου στο οποίο αναφέρεται το άρθρο.
  Λαμβάνεται από το περιεχόμενο της σελίδας βασιζόμενο ευρετικές μεθόδους όπως ο αριθμός των ανδρικών ή των θηλυκών αντωνυμών που χρησιμοποιούνται στο κείμενο.
- H. NUMBER, για την class wikibase-term Το αν η περιγραφόμενη έννοια αναφέρεται σε ενικό ή πληθυντικό αριθμό. Λαμβάνεται από το περιεχόμενο του κειμένου με βάση χαρακτηριστικά όπως το πόσες φορές ο τίτλος της σελίδας εμφανίζεται στον πληθυντικό ή στον ενικό αριθμό. Έχει αξία για όλα τα αντικείμενα. Επιστρέφει #t αν είναι πληθυντικός, #f αν είναι ενικός.
- I. PROPER, για την class wikibase-term
  Το αν η περιγραφόμενη έννοια είναι κύριο όνομα. Λαμβάνεται από το περιεχόμενο του κειμένου με βάση τα χαρακτηριστικά όπως το πόσες φορές ο τίτλος της σελίδας εμφανίζεται με κεφαλαία γράμματα όταν δεν είναι στην αρχή της σελίδας. Έχει τιμή για όλα τα αντικείμενα.
  Επιστρέφει #t αν είναι κύριο όνομα, #f αν δεν είναι.

#### 2. Η εντολή get-classes

Δεδομένου του ονόματος ενός αντιχειμένου, επιστρέφει μια λίστα με

όλες τις classes οπου ανήκει το αντικείμενο, με τις classes να αντιπροσωπεύονται ως lisp-readable strings. Παραδοσιακά τα ονόματα των classes δίνονται με μικρά γράμματα χωρίς όμως αυτό να είναι απολύτως απαραίτητο.

#### 3. Η εντολή get-attributes

Δεδομένου του ονόματος μιας class, επιστρέφει μια λίστα με όλα τα χαρακτηριστικά της class, ως lisp-readable strings. Τα ονόματα των χαρακτηριστικών δίνονται με κεφαλαία γράμματα, αλλά αυτό δεν αποτελεί απόλυτη απαίτηση.

#### 4. Η εντολή sort-symbols

Βάζει σε σειρά τα δεδομένα σύμβολα με βάση το μέγεθος του αντίστοιχου άρθρου, ομαδοποιώντας σύμβολα με ίσο μέγεθος άρθρου.

```
(sort-symbols "Obama<sub>□</sub>(surname)" "Barack<sub>□</sub>Obama")
=> (("Barack<sub>□</sub>Obama") ("Obama<sub>□</sub>(surname)"))
```

#### 5. Η εντολή sort-symbols-named

Παίρνει ένα σύμβολο  $\alpha$  και ένα σύνολο συμβόλων  $\beta_i$ . Βάζει τα  $\beta_i$  σε σειρά έτσι ώστε εάν κάποιο σύμβολο είναι το ίδιο με το  $\alpha$ , το ίδιο και το υποσύνολό του μπαίνουν στην αρχή.

```
(sort-symbols-named
```

```
"cake"
 "Cake_{\sqcup}(TV_{\sqcup}series)"
 "Cake_{\sqcup}(firework)"
 "Cake_(film)"
 "Cake_{\sqcup}(drug)"
 "Cake"
 "Cake_(band)"
 \texttt{"Cake}_{\sqcup}(\texttt{advertisement})\,\texttt{"}
 "The_{\sqcup}Cake")
=> (("Cake")
("Cake_(band)")
("Cake_{\sqcup}(advertisement)")
("Cake_{\sqcup}(TV_{\sqcup}series)")
("The_{\sqcup}Cake")
("Cake_(film)")
("Cake_{\sqcup}(firework)")
("Cake<sub>\(\)</sub>(drug)"))
```

### Chapter 2

# **Getting started**

Η συνολική WikipediaBase βρίσκεται σε ένα git repository στο infolab's github orginization page  $\{\{\text{ref(infolab}_{github})}\}\}$ .

```
git clone git@github.com:infolab-csail/WikipediaBase
```

Το WikipediaBase εξαφτάται από πολλά άλλα πακέτα python για τη λειτουφγία του. Ευτυχώς, η python είναι πακεταφισμένη όχι μόνο με ένα σημαντικό package manager (το pip) αλλά επίσης με ένα μηχανισμό που ονομάζεται virtualenv το οποίο απομονώνει την εγκατάσταση των εξαφτήσεων από το υπόλοιπο σύστημα. Έτσι αποφεύγονται πφοβλήματα όπως ασυμβατότητα εκδόσεων ή namespace collisions. Ο τφόπος που δουλεύει το virtualenv είναι αντιγφάφοντας ένα μέφος από το global python installation και κάνοντας symlink το υπόλοιπο σε ένα τοπικό φάκελο και εγκαθιστώντας τα dependencies στο τοπικό sandbox.

Ένα python virtualenv δημιουργείται και ενεργοποιείται ως εξής:

```
$ virtualenv --no-site-packages py
$ . py/bin/activate
$ which python
/the/local/directory/py/bin/python
```

Τώρα που ασφαλώς τα έχουμε εγκαταστήσει όλα θέλουμε χωρίς να σπάσουμε global installation

```
pip install -r requirements.txt
```

Θα χρειαστούμε μερικά επιπλέον εργαλεία για να δουλέψει η WikipediaBase που θα πρέπει να εγκατασταθούν system wide:

- Postresql
- Redis

Η εγκατάσταση αυτών των πακέτων διαφέρει ανάλογα με το λειτουργικό σύστημα ή τον package manager. Και οι δύο είναι βάσεις δεδομένων. Ο σκοπός τους είναι πρώτον, η προσωρινή αποθήκευση συχνά επαναλαμβανόμενων υπολογισμών (caching), και δεύτερον η αποθήκευση ahead-of-time υπολογισμών, όπως το START.

### Chapter 3

# Αρχιτεκτονική

Παρακάτω παρουσιάζονται τα μέρη του συστήματος WikipediaBase και ο τρόπος που αλληλεπιδρούν.

#### 1. Infobox

Για το σκοπό της παρούσας εργασίας θεωρούμε ένα infobox I με χαρακτηριστικά  $a_i$  και τιμές  $v_i$  είναι ένα σύνολο από ζεύγη  $a_i, v_i$  μαζί με ένα τύπο infobox t. Κάθε χαρακτηριστικό  $a_i$  και τιμή  $v_i$  έχουν 2 μορφές:

- rendered μορφή,  $a_i^r$  και  $v_i^r$  αντίστοιχα, η rendered HTML αναπαράσταση
- Η markup αναπαράσταση,  $a_i^m$  και  $v_i^m$  που είναι η mediawiki markup συμβολοσειρά

Η python class Infodox είναι ο βασικός τύπος δεδομένων για την πρόσβαση σε πληροφορίες από το infodox ενός άρθρου. Η Infodox, όπως και η Article, είναι αυτή που θα χρησιμοποιήσει κάποιος όταν χρησιμοποιεί τη wikipediabase ως βιβλιοθήκη Python. Οι μέθοδοι που παρέχονται από την Infodox δίνουν πρόσβαση στις εξής πληροφορίες:

Κλάσεις επειδή έχουμε δημιουργήσει python αντικείμενα Infobox βασισμένοι σε ένα όνομα συμβόλου (π.χ. όνομα της σελίδας) το οποίο στο άρθρο του μπορεί να έχει παραπάνω από ένα wikipedia infoboxes διαφορετικών κλάσεων, ένα python αντικείμενο Infobox μπορεί στην πραγματικότητα να είναι μια διεπαφή για πολλαπλά wikipedia infoboxes. Για την ανάκτηση μιας symbol class σε μορφή κατάλληλη για το START, υπάρχει μια διαφορετική μέθοδος.

Τιμές χαρακτηριστικών δηλαδή είτε  $v_i^r$  είτε  $v_i^m$  δεδομένου είτε  $a_i^r$  είτε  $a_i^m$ .

**Ονόματα χαρακτηριστικών** που παρέχονται με τη χρήση του MetaInfobox (βλέπε παρακάτω)

Εξαγωγή των πληφοφοφιών σε python types συγκεκφιμένα

- dict yıa  $a_i^r \rightarrow v_i^r \circ a_i^m \rightarrow v_i^m$
- Το συνολικό infobox rendered, ή σε ένα markup μορφή.

Τα Infoboxes οργανώνονται σε μια ευρεία ιεραρχία η οποία στον κώδικα του WikiepdiaBase αναφέρεται ως infobox tree. Το infobox tree ανακτάται από σελίδα wikipedia List of infoboxes και χρησιμοποιείται για να συνταχθεί η οντολογία των όρων wikipedia δηλαδή η κατάταξή τους σε κλάσεις.

#### 2. MetaInfobox

Το MetaInfodox υλοποιείται ως μια υποκλάσση του Infodox και προσδίδει πληροφορία σχετικά με το infodox, εστιάζοντας στην αντιστοιχία της rendered μορφής των χαρακτηριστικών με την αντίστοιχη markup μορφή. Έτσι δεδομένου ενός infodox τύπου I έχει πιθανά χαρακτηριστικά  $a_1,...,a_n$ . Κάθε χαρακτηριστικό έχει δύο αναπαραστάσεις:

- τη markup αναπαράσταση που χρησιμοποιείται στο infobox template.
- την HTML rendered αναπαράσταση, που είναι το κείμενο που φαίνεται στην αριστερή μεριά του πίνακα του infobox στη σελίδα της wikipedia.

Παραδείγματος χάριν στα officeholder infoboxes υπάρχει ένα χαρακτηριστικό με markup αναπαράσταση predecessor που έχει rendered αναπαράσταση Preceded by.

Για να το πετύχει αυτό το MetaInfodox βρίσκει το markup representation όλων των αποδεκτών χαρακτηριστικών μιας κλάσης infodox μέσα από την σελίδα τεκμηρίωσης του αντίστοιχου template. Στη συνέχεια δημιουργεί ένα infodox όπου κάθε χαρακτηριστικό έχει ως τιμή τη markup αναπαράσταση του χαρακτηριστικού αυτού, προσθέτοντας πριν και μετά τη συμβολοσειρά!!!. (Για παράδειγμα το χαρακτηριστικό με markup όνομα predecessor θα έχει τιμή!!!predecessor!!!). Στη συνέχεια κάνει render το infodox που δημιούργησε και ψάχνει για τη συμβολοσειρά!!!predecessor!!! στις rendered τιμές. Θεωρούμε ότι οι τα αντίστοιχα rendered ονόματα αντιστοιχούν στα markup χαρακτηριστικά. Πρέπει να σημειωθεί ότι η αντιστοιχία των rendered χαρακτηριστικών με τα markup χαρακτηριστικά δεν είναι αμφοσήμαντη, δηλαδή κάθε

markup χαρακτηριστικό μπορεί να αντιστοιχεί σε μηδέν η περισσότερα rendered χαρακτηριστικά και το αντίστροφο.

Για παράδειγμα για ένα infobox τύπου Foo με αποδεκτά χαρακτηριστικά A,B,C και D το MetaInfobox θα δημιουργούσε markup:

```
{{Infobox Foo
| A = !!!A!!!
| B = !!!B!!!
| C = !!!C!!!
| D = !!!D!!!
```

Και η rendered μορφή θα ήταν, ανάλογα με την υλοποίηση του Foo infobox.

Attribute	Value
A	!!!A!!! !!!B!!! !!!C!!!
В	!!!A!!! !!!B!!! !!!C!!!
C	!!!A!!! !!!B!!! !!!C!!!
D	!!!D!!!

Έτσι η αντιστοιχία γίνεται σχετικά εμφανής.

#### 3. Article

Η class Article είναι υπεύθυνη για την πρόσβαση σε κάθε πληροφορία σχετική με το άρθρο γενικότερα. Αυτό περιλαμβάνει τις παραγράφους, επικεφαλίδες, τον πηγαίο markup κωδικα και τις κατηγορίες MediaWiki.

#### 4. Fetcher

Η κλάση Fetcher αναλαμβάνει την επικοινωνία της WikipediaBase με τις πηγές πληφοφοφιών. Είναι ένα μονήφες αντικείμενο που υλοποιεί μια συγκεκφιμένη διεπαφή.

Τα υλοποιημένα Fetchers έχουν μια κληφονομική ιεφαρχία που φαίνεται από την παρακάτω λίστα.

BaseFetcher είναι η υπερκλάση όλων των fetchers. Θα επιστρέψει αυτούσιο το symbol. Κάνουμε override αυτή τη λειτουργία στις κληρονόμους κλάσεις για να υλοποιήσουμε τη λογική της διεπαφής με τις πηγές πληροφοριών.

Fetcher Υλοποιεί τη βασική λειτουργία. Αναζητά πληροφορίες απο το wikipedia.org. Είναι δυνατόν να κατευθύνουμε ένα Fetcher προς ένα wikipedia mirror αλλά η εκτέλεση σε wikipedia-mirror είναι απαγορευτική από άποψη επίδοσης.

Caching Fetcher κληφονομεί από την class Fetcher και διατηφεί τη λειτουργικότητα, μόνο που χρησιμοποιεί μια βάση δεδομένων για την προσωρινή αποθήκευση των πληφοφοριών. Είναι η προεπιλεγμένη fetcher class.

StaticFetcher είναι μια κλάση που υλοποιεί το interface της BaseFetcher αλλά αντί να φτάσει σε κάποια πηγή πληφοφοφιών για τα δεδομένα δίνει τιμές επιστφοφής στατικά οφισμένες. Χφησιμοποιείται κυφίως από το MetaInfobox.

Από ποοεπιλογή, το markup ποοέοχεται από μια βάση δεδομένων. Αν η παράμετρος force\_live του constructor έχει οριστεί σε True τότε το markup θα ληφθεί από το wikipedia.org. Όταν οι δοκιμές τρέχουν στο TravisCI[27], θέλουμε πάντα να χρησιμοποιούνται live δεδομένα. Ελέγχουμε αν ο Travis εκτελεί δοκιμές κοιτάζοντας τη μεταβλητή περιβάλλοντος WIKIPEDIABASE\_FORCE\_LIVE.

#### 5. Renderer

Οι Renderers είναι μονήσεις classes, χρήσιμες για την μετατροπή MediaWiki markup σε HTML. Αρχικά χρησιμοποιήθηκε για την μετατροπή το wikiepedia sandbox[31], επειδή είναι ελαφρώς ταχύτερο από το Wikipedia API. Μεταπηδήσαμε στο wikipedia.org API γιατί το wikipedia-mirror ήταν πολύ αργό και το wikipedia.org θεώρησε κατάχρηση της υπηρεσίας με αποτέλεσμα να μπλοκάρει το IP μας μετά από μερικά τεστ. Γι' αυτό το λόγο χρησιμοποιήθηκε τελικά το API, με Redis caching. Αυτό λειτούργησε αρκετά καλά, επειδή τα Renderer αντικείμενα καταλήγουν να χρησιμοποιούνται μόνο από το MetaInfobox, το οποίο έχει ένα αρκετά περιορισμένο πεδίο εφαρμογής, και έτσι τα cache misses είναι σπάνια.

Μια ενδιαφέρουσα πληφοφορία για την class Renderer ήταν ότι αυτός ήταν ο λόγος που ένα ζευγάρι CSAIL αποκλείστηκε προσωρινά από την επεξεργασία της wikipedia. Ενώ η wikipedia.org έχει μια πολύ επιεική πολιτική όταν πρόκειται για τον αποκλεισμό των χρηστών που έχουν κάνει spamming τους servers, επαναλαμβανόμενες δοκιμές της κατηγορίας Renderer με στόχευση το wikipedia sandbox προκάλεσε το IP του δοκιμαστικού μηχανήματος να αποκλεισθεί προσωρινά με το σκεπτικό ότι "η δραστηριότητα του δεν προάγει την βελτίωση της wikipedia". Επανατοποθετήσαμε το

Renderer να χοησιμοποιεί το wikipedia API και ποτέ δεν είχαμε ξανά πρόβλημα με την ρύθμιση της wikipedia.

#### 6. Pipeline

Κατά την επίλυση ενός ερωτήματος η WikipediaBase ενεργοποιεί ένα pipeline λειτουργιών για να διαπιστωθεί ποιος είναι ο καλύτερος τρόπος απάντησης.

#### (a) Frontend

Η WikipediaBase μπορεί να χρησιμοποιηθεί ως βιβλιοθήκη αλλά η πρωταρχική της λειτουργία είναι ως backend στο START. Η επικοινωνία μεταξύ START και WikipediaBase γίνεται πάνω από μια plaintext telnet σύνδεσή στην πόρτα 8023 χρησιμοποιώντας s-expressions. Το frontend χειρίζεται το δίκτυο σύνδεσης με το START, μεταφράζει τις προσλαμβανόμενες ερωτήσεις σε κλήσεις της Knowledgebase και στη συνέχεια μεταφράζει την αντίδραση της Knowledgebase σε κατάλληλα διαμορφωμένες εκφράσεις και τις επιστρέφει πίσω στο telnet connection.

#### (b) Knowledgebase

Η knowledgebase είναι το σημείο εισαγωγής στην υπόλοιπη wikipediahase

Χοησιμοποιεί μοτίβο Provider/Acquirer (βλ. παρακάτω) για να παρέχει διαφανή διεπαφή της frontend με αυθαίρετες μεθόδους. Οι μέθοδοι αυτοί είναι υπεύθυνοι για την επιλογή του αν θέλουμε να καταλήξουμε σε classifiers, resolvers ή οποιοδήποτε άλλο μηχανισμό για να δοθεί απάντηση στο ερώτημα που τέθηκε. Οι διαθέσιμοι classifiers και resolvers γίνονται προσβάσιμοι αυτόματα στη knowledgebase χρησιμοποιώντας τη βασική τους κλάση.

#### (c) Classifiers

Κάθε Classifier είναι μονήφης κλάση και υλοποιεί μια ευφετική μέθοδο για να συντάξει μια λίστα από classes ενός symbol. Ένα symbol μποφεί να επιστφέφει μηδέν ή πεφισσότεφες classes. Συνήθως, ένας Classifier θα επιλέξει μόνο αν ένα αντικείμενο πφάγματι ανήκει σε μια συγκεκφιμένη κατηγοφία ή όχι, αλλά αυτό δεν είναι απαφαίτητο.

#### i. Term

Ο TermClassifier απλά αναθέτει την κατηγορία wikipedia-term. Η Wikipediabase διαπραγματεύεται μόνο με πληροφορίες σχετικές με τη wikipedia. Συνεπώς όλες οι έννοιες που συναντώνται ανήκουν σε αυτήν την κατηγορία.

#### ii. Infobox

Το InfoboxClassifier αναθέτει σε ένα symbol την κατηγορία infobox. Για παράδειγμα η σελίδα Bill Clinton περιέχει το infobox:

Και γι αυτό λαμβάνει την κατηγορία wikipedia-president.

#### iii. Person

To PersonClassifier αναθέτει την κατηγορία wikibase-person χρησιμοποιώντας κάποια χαρκτηριστικά με την σειρά που περιγράφονται:

- Category regex matches
- Category regex excludes
- · Category matches

Περιγράφονται λεπτομερώς στο παράρτημα.

#### (d) Resolvers

Οι Resolvers είναι επίσης μονήσεις κλάσεις αλλά ο σκοπός τους είναι να βρούν την τιμή του αναζητούμενου χαρακτηριστικού. Όλοι οι resolvers κληφονομούν από την class BaseResolver και πρέπει να υλοποιούν τις ακόλουθες μεθόδους:

- resolve(class, symbol, attribute) που δίνει την τιμή ενός χαρακτηριστικού δεδομένου του symbol και της class.
- attributes(class, symbol): που δίνει μια λίστα από τα χαρακτηριστικά που μπορεί να επιλύσει ο συγκεκριμένος resolver για το συγκεκριμένο άρθρο δεδομένης της class του.

Οι υλοποιημένοι resolvers είναι οι απόλουθοι:

Error ο ελάχιστης προτεφαιότητας resolver. Επιλύεται πάντα σε σφάλμα.

Infobox Επιλύει χαρακτηριστικά που αναφέρονται σε κάποιο πεδίο του infobox

Person επιλύει τα ακόλουθα ειδικά χαρακτηριστικά των άρθρων που αναφέρονται σε πρόσωπα

- birth-date
- death-date
- gender

Sections το περιεχόμενοτων κεφαλαίων σε ένα άρθρο.

Term επιλύει ένα συγκεκριμένο σύνολο χαρακτηριστικών,

- coordinates Οι συντεταγμένες μιας γεωγραφικής περιοχής
- image Την εικόνα μέσα στο infobox.
- number Αληθής τιμή αν το σύμβολο είναι στον πληθυντικό (πχ The Beatles)
- proper Αληθής αν αναφέρεται σε κύριο όνομα.
- short-article Περίληψη του άρθοου, τυπικά η πρώτη παράγραφος.
- url  $H \delta \iota \epsilon \dot{\nu} \theta v v \sigma \eta \tau o v \dot{\alpha} \varrho \theta \varrho o v$ .
- word-cout Το μέγεθος του άρθρου σε λέξεις.

#### 7. Lisp types

Ο τύπος Lisp είναι περιτυλίγματα (wrappers) για python αντιχείμενα ή τιμές που παρουσιάζονται σε μορφή s-expression που το START μπορεί να κατανοήσει. Έχουν δημιουργηθεί είτε από το ανεπεξέργαστο ερώτημα και έχουν ξετυλιχθεί (unwrapped) ώστε να είναι χρήσιμα στο pipeline, ή από την απάντηση που δίνει η WikipediaBase και στη συνέχεια κωδικοποιούνται σε ένα string και αποστέλλονται μέσω telnet στο START.

## **Chapter 4**

# Το μοντέλο provider/acqirer

Η WikipediaBase προσπαθεί να είναι modular και με δυνατότητα επέκτασης. Για να επιτευχθεί αυτό, συχνά είναι χρήσιμο να συμπλέκει πολλαπλές πηγές του ίδιου τύπου δεδομένων. Αυτό είναι ιδιαίτερα χρήσιμο κατά την πρόσβαση ευρετικών μεθόδων όπως των classifiers που είδαμε παραπάνω. Για την προώθηση του modularity και για να αποφευχθεί ισχυρή αλληλεξάρτηση των υποσυστημάτων δημιουργήθηκε το μοντέλο provider/acquirer.

Ο Provider είναι ένα αντιπείμενο μέσω του οποίου μπορούμε να διαχειριστούμε πηγές που είναι αποθηπευμένες ως ζεύγη πλειδιού - τιμής. Η class Provider προσφέρει python decorators για να πάνει αυτή τη διάταξη εύπολη για τον προγραμματιστή. Ένας Acquirer έχει διαφανή (transparent) πρόσβαση στους πόρους πολλαπλών Providers σαν να ήταν ένα ενιαίο σύνολο πλειδιών. Αυτό το πρότυπο πυρίως χρησιμοποιείται για την KnowledgeBase ώστε να παρέχει στο Frontend ενιαίο τρόπο πρόσβασης στις πηγές.

#### 1. Παράδειγμα

Εκθέτουμε το μοτίβο provider/acquirer με ένα παράδειγμα ενθέτοντας μια μικρή lisp μέσα στην python, και χειριζόμενοι το state του εκτελούμενου προγράμματος με providers και acquirers.

```
from wikipediabase.provider import Provider,
    Acquirer, provide

class EvalContext(Acquirer):
    def __init__(self, closures):
        super(EvalContext, self).__init__(closures)
        self.closures = closures
```

```
def __call__(self, _ctx, expr):
        if isinstance(expr, list):
            # Handle quotes
            if expr[0] is 'quote':
                return expr[1]
            # Call the lambda
            fn = self(_ctx, expr[0])
            return fn(self, *[self(_ctx, e) for e in
                expr[1:]])
        if isinstance(expr, basestring) and expr in
           self.resources():
            return self(_ctx, self.resources()[expr
               ])
        return expr
class Lambda(Acquirer):
    def __init__(self, args, expr, env):
        # Get your symbols from all the available
           closures plus an
        # extra for local variables
        super(Lambda, self).__init__([env] + [
           Symbols()])
        self.args = args
        self.expr = expr
    def __call__(self, _ctx, *args):
        \# Add another closure to the list
        arg_provider = Provider();
        for s, v in zip(self.args, args):
            arg_provider.provide(s, v)
        # Build an eval context and run it
        ctx = EvalContext([arg_provider, Provider(
           self.resources())])
        return [ctx(ctx, e) for e in self.expr][-1]
class Symbols(Provider):
    @provide('setq')
    def setq(self, ctx, symbol, val):
        self.provide(symbol, val)
```

Αυτή η μικοή lisp αν και πρωτόγονη υποστηρίζει:

- lambdas
- A global symbol table
- · lexical scoping
- · conditionals
- Quoted literals

Ποοφανώς δεν είναι μια χρήσιμη γλώσσα αλλά μπορεί να πετύχει μερικά ενδιαφέροντα κόλπα:

Μπορούμε να χρησιμοποιήσουμε python types:

```
>>> GLOBAL_EVAL({}, 1)
1
>>> GLOBAL_EVAL({}, True)
True
>>> GLOBAL_EVAL({}, "hello")
'hello'
>>> GLOBAL_EVAL({}, list)
<type 'list'>
```

Μπορούμε να ορίσουμε lambdas και να τις καλέσουμε. Το ακόλουθο παράδειγμα είναι ισοδύναμο με το  $(\lambda a.a)1$ , το οποίο πρέπει να εκτιμηθεί στην τιμή 1:

```
>>> GLOBAL_EVAL({}, [["lambda", ['quote', ['a']], 'a
     '], 1])
1
```

Η μικοή μας lisp δεν είναι pure εφ όσον έχουμε mutable global symbol table. Αυτό σημαίνει πως η σειρά των διεργασιών έχει σημασία. Εφ όσον δεν έχουμε progn η άλλα macros συνηθισμένα σε lisp dialects ο καλύτερος τρόπος να κάνουμε διεργασίες σε σειρά είναι να τις εντάξουμε σε ένα lambda και να το εκτιμήσουμε (evaluate).

```
>>> GLOBAL_EVAL({}, [['lambda', ['quote', []], ['setq', 'b', 2], 'b']])
2
```

Ο προσεκτικός αναγνώστης ίσως παρατηρήσει ότι η λίστα για τα lambda arguments είναι quoted. Ο λόγος γι αυτό είναι ότι δεν θέλουμε η λίστα να εκτιμηθεί.

Συνεχίζοντας την έμθεση του provider/acqirer. Σε μάθε σημείο του μώδιμα το μάθε σύμβολο λαμβάνει τιμές από πολλαπλές πηγές. Με σειφά προτεφαιότητας:

- · The local closure
- The arguments of the lambda
- · Builtin functions

Όλα τα προηγούμενα εκτίθενται περιληπτικά χρησιμοποιώντας το provideraquirer model.

Σε κάθε σημείο ένα διαφορετικό EvaluationContext είναι υπεύθυνο για την εκτίμηση και κάθε EvaluationContext έχει πρόσβαση στα γνωστά σύμβολα του μέσω μιας array of providers τα οποία εκτίθενται περιληπτικά χρησιμοποιώντας το υπό συζήτηση μοντέλο.

### Chapter 5

# **Testing**

Η καλή λειτουργία της WikipediaBase εξασφαλίζεται από μια ολοκληφωμένη σειρά δοκιμών, των unit tests, functional tests και regression tests. Τα Unit tests ελέγχουν μια μικρή ομάδα του functionality, το οποίο έχει συντεθεί για την δημιουργία του όλου συστήματος της WikipediaBase. Για το unit testing χρησιμοποιούμε την default βιβλιοθήκη python για testing. Κάθε τεστ είναι μια class μου κληρονομεί από την class TestCase και υλοποιεί το interface της που περιγράφεται παρακάτω.

Τα Functional tests είναι γραμμένα από ποιν, κατά τη διάρκεια ή λίγο μετά τη δημιουργία του συστήματος που τεστάρουν και επιβεβαιώνουν τη σωστή συνολική λειτουργία του συστήματος. Τα Regression tests είναι πολύ παρόμοια με τα to functional tests. Αποδεικνύουν ότι όταν βρεθεί ένα σφάλμα (bug) αυτό διορθώθηκε και επιβεβαιώνουν ότι δεν θα εμφανισθεί ξανά αργότερα. Τα Functional και τα regression tests είναι τοποθετημένα στα tests/examples.py Σχεδόν όλα τα τεστ ξεκινούν με τον ακόλουθο κώδικα:

```
from __future__ import unicode_literals
try:
    import unittest2 as unittest
except ImportError:
    import unittest
from wikipediabase import fetcher
```

Το παραπάνω είναι ειδικό για το the fetcher module. Όπως είναι προφανές χρησιμοποιούμε το unittest module από την βιβλιοθήκη python. Το test το ίδιο έχει το ακόλουθο format:

```
class TestFetcher(unittest.TestCase):
    def setUp(self):
        self.fetcher = fetcher.get_fetcher()

    def test_html(self):
        html = self.fetcher.html_source("Led_Zeppelin")
        self.assertIn("Jimmy_Page", html)
```

Η setUp μέθοδος τρέχει πριν από κάθε τεστ του TestCase. Τα τεστ του testcase αντιπροσωπεύονται από μεθόδους της class των οποίων το όνομα αρχίζει με test\... Στην συγκεκριμένη περίπτωση παίρνουμε την σελίδα της wikipedia για το Led Zeppelin και τσεκάρουμε ότι το όνομα Jimmy Page αναφέρεται τουλάχιστον μια φορά. Αυτό φανερά δεν αποδεικνύει ότι το fetcher δεν φέρνει για παράδειγμα την σελίδα για το Yardbirds, Page's first band. Για αυτό το λόγο γράφουμε παραπάνω από ένα αυτού του είδους τεστ.

Στην περίπτωση του fetcher, για να ακολουθήσουμε το παραπάνω παράδειγμα, το συνολικό τεστ υπάρχει στο παράρτημα.

Εφαρμόσαμε το εργαλείο nosetests να βρούμε και να τρέξουμε τα τεστ. Για να το κάνουμε αυτό το προσθέσαμε σαν προαπαιτούμενο στο setup.py.

```
from setuptools import setup

setup(
    tests_require=[
        'nose>=1.0',
        ...
    ],
    ...
    test_suite='nose.collector',
    ...
)

Στη συνέχει να τφέξουμε τα τεστ:

$ python setup.py test
```

Η Nose θα βοει όλα τα αρχεία τα οποία είναι στο φάκελο tests/ και έχουν το πρόθεμα test\\_, για παράδειγμα test\\_fetcher.py. Μέσα σ αυτά τα αρχεία η nose θα αναζητήσει subclass της TestCase και των οποίων το όνομα

αρχίζει με Test, για παράδειγμα TestFetcher. Στη συνέχεια τρέχει όλες τις μεθόδους από τις collected classes που έχουν το πρόθεμα test\\_. Είναι επίσης δυνατό να τρέξει μόνο συγκεκριμένα τεστ.

```
$ python setup.py test --help
Common commands: (see '--help-commands' for more)
 setup.py build will build the package underneath
     'build/'
 setup.py install will install the package
Global options:
 --verbose (-v) run verbosely (default)
 --quiet (-q) run quietly (turns verbosity off)
 --dry-run (-n) don't actually do anything
 --help (-h) show detailed help message
 --no-user-cfg ignore pydistutils.cfg in your home
     directory
Options for 'test' command:
 --test-module (-m) Run 'test_suite' in specified
     module
  --test-suite (-s) Test suite to run (e.g. '
     some_module.test_suite')
 --test-runner (-r) Test runner to use
usage: setup.py [global_opts] cmd1 [cmd1_opts] [cmd2 [
   cmd2_opts] ...]
  or: setup.py --help [cmd1 cmd2 ...]
  or: setup.py --help-commands
  or: setup.py cmd --help
```

Δείτε το παράρτημα για επιτυχημένη εκτέλεση των τεστ.

### Chapter 6

# Συνώνυμα

Ποιν μιλήσουμε για τα συνώνυμα είναι σημαντικό να ορίσουμε πιο αυστηρά τα symbols στο πεδίο του omnibase universe:

Σύμβολα είναι ταυτοποιητές των "αντικειμένων", "objects", στις πηγές των πληφοφοφιών (ο όφος "σύμβολο" ("symbol") είναι ατυχής γιατί έχει διάφοφες έννοιες στην επιστήμη των υπολογιστών. Δυστυχώς έχει μείνει για ιστοφικούς λόγους.)

Δεδομένου ότι η γλώσσα τείνει να έχει πολλαπλές λέξεις που αναφέφονται στο ίδιο πράγμα, είναι επιτακτική η ανάγκη να καθορισθούν πολλά ονόματα για κάθε σύμβολο. Συνώνυμα είναι τα ονόματα τα οποία οι χρήστες μπορούν να χρησιμοποιήσουν για να αναφερθούν σε ένα συγκεκριμένο σύμβολο.

(Ο όφος συνώνυμα "synonym" είναι ατυχής γιατί είναι one-way mapping - "gloss" θα ήταν καλύτεφος όφος αλλά έμεινε ο όφος συνώνυμα για ιστοφικούς λόγους)

Ο ορισμός συνωνύμων είναι δουλειά του backend. Για το λόγο αυτό αναλαμβάνει η WikipediaBase να ορίσει τα απαιτούμενα συνώνυμα.

#### 1. Καλά και κακά συνώνυμα

Υπάρχουν κανόνες για το ποιο είναι καλό ή κακό συνώνυμο

- Δεν πρέπει να ξεκινούν με άρθρα ("the", "a", "an")
- Δεν πρέπει να ξεκινούν με "File:" or "TimedText:".
- Δεν πρέπει να περιέχουν HTML anchors. Πχ "Alexander<sub>Pushkin</sub>#Legacy"
- Δεν πρέπει να ξεκινούν με τα ακόλουθα:
  - "List of "
  - "Lists of"

- "Wikipedia: "
- "Category: "
- ":Category: "
- "User: "
- "Image: "
- "Media: "
- "Arbitration in location"
- "Communications in location"
- "Constitutional history of location"
- "Economy of location"
- "Demographics of location"
- "Foreign relations of location"
- "Geography of location"
- "History of location"
- "Military of location"
- "Politics of location"
- "Transport in location"
- "Outline of topic"
- Δεν πρέπει να ταιριάζει \d\d\d in location ή location in \d\d\d
- Δεν πρέπει να είναι ονόματα των disabiguation pages. Για να το κάνουμε αυτό έτσι ώστε να συμπεριλαμβάνει όλες τις σχετικές σελίδες (συμπεριλαμβανομένων των τυπογραφικών λαθών) εννοούμε σύμβολα που ταιριάζουν με \([Dd]isambig[\^)]\*\).
- Συνώνυμα που α) θα μπορούσαν να εκληφθούν ότι ξεκινούν με άρθρα και β) μπορεί να εκλείπουν κάτι χρήσιμο. Αυτό σημαίνει ότι για παράδειγμα "Α. House" (συνώνυμο του «Abraham House") είναι ελλιπών προδιαγραφών διότι ενδέχεται να παραπλανήσει το START στην περίπτωση των ερωτήσεων όπως "Πόσο κοστίζει ένα σπίτι στη Silicon Valley;". Αφετέρου "a priori" μπορεί να διατηρηθεί επειδή δεν υπάρχουν λογικές ερωτήματα όπου "α" είναι ένα άρθρο πριν "priori".

#### 2. Παραγωγή συνωνύμων

Για να συμβιβάσουμε αυτούς τους περιορισμούς δύο μέθοδοι χρησιμοποιούνται: qualification και modification των υποψήφιων συνωνύμων. Πρώτα προσπαθούμε τη modification μέθοδο και αν αυτό αποτύχει επιχειρούμε να εκτελέσουμε qualification. Οι κανόνες για modification έχουν ως εξής:

- Να διαγράψουμε τα άρθρα από την αρχή ενός συνωνύμου:
  - "A "
  - "An "
  - "The"
  - "(The)"
  - The
  - κτλ
- Να δημιουργούμε και και τα δύο versions, με και χωρίς παρενθέσεις.
   Πχ, δεδομένου του συμβόλου "Raven (journal)" δημιουργούμε:
  - "Raven (journal)"
  - "Raven"
- Να χοησιμοποιήσουμε τη συμβολοσειρά ποιν και μετά το slash, αλλά όχι το αρχικό symbol, πχ. δεδομένου του symbol "Russian language/Russian alphabet" δημιουργούμε
  - "Russian language"
  - "Russian alphabet"
- Να αναστρέψουμε των ανεστραμμένων συμβόλων με κόμματα. Πχ δεδομένου "Congo, Democratic Republic Of The", αναστρέφουμε για να πάρουμε "Democratic Republic Of The Congo"
- Ώς συνήθως, να απορρίψουμε leading articles εάν είναι αναγκαίο. Π.χ. δοθέντος συμβόλου "Golden ratio, the" το αντικαθιστούμε με "the Golden ratio", στη συνέχεια διαγράφουμε τα άρα για πάρουμε: "Golden ratio" το ίδιο συμβάνει για τα a, an, κτλ.

Με αυτό τον τφόπο δημιουργήσαμε ένα αρχικό πακέτο συνωνύμων από το ίδιο το όνομα του αντικειμένου. Επιπλέον μποφούμε να δημιουργήσουμε άλλο ένα πακέτο από τα wikipedia redirects στο άφθφο. Η Wikipedia παφέχει ένα SQL dump για όλα τα redirects. Για να φοφτώσουμε τον πίνακα του SQL dump σε μια βάση δεδομένων όπου έχουμε φοφτώσει ήδη τα δεδομένα της wikipedia:

```
wget https://dumps.wikimedia.org/enwiki/latest/
   enwiki-latest-redirect.sql.gz \
-0 redirect.sql.gz && gzcat redirect.sql.gz |
    mysql
```

Και στη συνέχεια μπορούμε να τρέξομε το SQL query για να βρούμε όλα τα συνώνυμα του (καλά και κακά) Bill Clinton:

```
select page_title, rd_title from redirect join page
  on rd_from = page_id and (rd_title = "
   Bill_Clinton" or page_title = "Bill_Clinton");
```

Για το πλήφες output δείτε στο παφάφτημα.

### **Databases and data sources**

#### 1. HTML and MediaWiki API

Η αρχική προσέγγιση για να πάρουμε τα δεδομένα της wikipedia είναι να ανασύρουμε τις φυσιολογικές HTML εκδόσεις των άρθρων της wikipedia και χρησιμοποιώντας edit pages να ανασύρουμε το mediawiki markup. Αρχικά χρησιμοποιήσαμε το αρχικό wikipedia.org site για λόγους performance (Βλέπε κεφάλαιο wikipedia-mirror runtime performance).

Το Mediawiki παρέχει ένα RESTful API για όλη την απαιτούμενη λειτουργία της wikipedia. Η βασική αρχή είναι ότι κάποιος μπορεί να στείλει αιτήματα με μεθόδους POST ή GET και να λαμβάνει απάντηση με την μορφή XML ή JSON. Η προτιμητέα απάντηση για την WikipediaBase ήταν να στέλνουμε GET HTTP αιτήματα και να λαμβάνουμε JSON δεδομένα. Το GET επιλέχθηκε επειδή προτάθηκε στην mediawiki API page γιατί το caching συμβαίνει στο HTTP επίπεδο. Σύμφωνα με τις οδηγίες του HTTP τα POST αιτήματα δεν μπορούν να είναι cached. Για το λόγο αυτό όταν διαβάζει κάποιος δεδομένα από web service API, θα πρέπει να χρησιμοποιεί GET αιτήματα και όχι POST.

Επίσης πφέπει να σημειωθεί ότι ένα αίτημα δεν μποφεί να εκτελεσθεί από cache εκτός αν το URL είναι ακφιβώς το ίδιο. Εάν ζητήσει κάνεις ένα αίτημα για api.php?titles=Foo|Bar|Hello, και αποθηκεύσει το αποτέλεσμα, μετά api.php?titles=Hello|Bar|Hello|Foo δεν θα βφει την απάντηση στην cache παφ όλο που είναι το ίδιο αιτήμα!

Η αναπαράσταση JSON επιλέχθηκε άπλα επειδή η βιβλιοθήκη json της python πολύ πιο εύκολη στη χρήση από την lxml, τη βιβλιοθήκη που χρησιμοποιούμε για XML/HTML parsing.

#### 2. Caching

Η Wikipediabase χρησιμοποιεί κυρίως έναν απομακρυσμένο χώρο αποθήκευσης δεδομένων που εφαρμόζει το mediawiki interface (δηλαδή το mediawiki). Προσπαθεί να αντιμετωπίσει ζητήματα επιδόσεων που προκύπτουν με την προσωρινή αποθήκευση των σελίδων σε μια τοπική key-value βάση δεδομένων. Το interface με τη βάση δεδομένων αφαιρείται με τη χρήση ενός python dictionary-style interface, το οποίο υλοποιείται στο persistentky.py. Ένα άλλο χαρακτηριστικό που το interface στην βάση δεδομένων πρέπει να υλοποιεί είναι η κωδικοποίηση των αποθηκευμένων αντικειμένων. Επειδή όλη η αποθηκευμένη πηροφορία είναι κείμενο, η βάση δεδομένων πρέπει να είναι ικανή να ανασύρει ακριβώς το κείμενο που έχει αποθηκευθεί λαμβάνοντας υπόψη την κωδικοποίηση. Λόγω των περιορισμών του DBM's τα κλειδιά (keys) πρέπει να είναι μόνο κωδικοποιημένα ASCII. Η βασική class για αλληλεπίδραση με την βάση δεδομένων, το EncodedDict, εφρμόζει τις μεθόδους \_encode\_key και \_decode\_key.

#### (a) DBM

Διάφορες υλοποιήσεις dbm[5] παρέχονται από την standard βιβλιοθήμη της python. Μεριμές διαθέσιμες εφαρμογές DBM είναι:

- AnyDBM
- GNU DBM
- Berkeley DBM

Είναι σημαντικό να αναφέφουμε ότι η ομαλή λειτουργία αυτών των βιβλιοθηκών εξαρτάται σε σημαντικό βαθμό από την βασική πλατφόρμα όπως το λειτουργικό. Όπως αναφέρθηκε παραπάνω οι interface classes του DBM μεταφράζουν από και προς ASCII.

#### (b) SQLite

Η SQLite[24] επίσης χρησιμοποιείται ως caching backend βάση δεδομένων. Δυστυχώς η αποτελεσματικότητά του στο δικό μας σκοπό ήταν απογοητευτική. Χρησιμοποιήσαμε ένα πολύ λεπτό wrapper, το sqlitedict[25], για να πάρουμε ένα key-value interface στην SQLite – μια relational βάση δεδομένων. Ο σχετικός WikipediaBase κώδικας είναι πολύ σύντομος:

Παρακάτω είναι δυο benchmark functions που θα διαβάσουν και θα γράψουν 100000 φορές στην βάση.

```
def benchmark_write(dic, times=100000):
    for i in xrange(times):
        dic['o' + str(i)] = str(i) * 1000

def benchmark_read(dic, times=100000):
    for i in xrange(times):
dic['o' + str(i)]
```

Και παρακάτω φαίνεται πως συγκρίνονται τα διάφορα backends χρησιμοποιώντας αυτές τις δυο συναρτήσεις.

```
>>> import timeit
>>> sqlkv = SqlitePersistentDict('/tmp/bench1.
   sqlite')
>>> timeit.timeit(lambda : benchmark_write(sqlkv
   ), number=100)
10.847157955169678
>>> timeit.timeit(lambda : benchmark_read(sqlkv)
   , number=100)
18.88098978996277
>>> dbmkv = DbmPersistentDict('/tmp/bench.dbm')
>>> timeit.timeit(lambda : benchmark_write(dbmkv
   ), number=100)
0.18030309677124023
>>> timeit.timeit(lambda : benchmark_read(dbmkv)
   , number=100)
0.14914202690124512
```

Η DBM βαση δεδομένων είναι σχεδόν 100 φορές ταχύτερη από

sqlite. Η διαφορά στην εκτέλεση οφείλεται στις διαφορετικές committing policies που έχουν μεταξύ τους. Μπορεί να είναι δυνατόν να ρυθμιστεί το SQLite ώστε να είναι τόσο γρήγορο όσο η DBM αλλά όχι με κάποιον εύκολο τρόπο.

#### (c) Άλλα backends

Και άλλα backends λαμβάνονται υπόψη, πυρίως το Redis το οποίο εφαρμόσθηκε αμέσως μετά την παράδοση της εργασίας από τον Alvaro Morales. Ο λόγος που αρχικά δεν το χρησιμοποιήσαμε ήταν γιατί έχει μοντελοποιηθεί ως ένας server-client και προσθέτει περιπλοκότητα σε ένα τμήμα του συστήματος το οποίο πρέπει να είναι όσο το δυνατόν πιο απλό. Ένας άλλος λόγος του αρχικού προβληματισμού μας ήταν σχετικά με το ότι το redis είναι ανεξάρτητο project δηλαδή δεν είναι μέρος της python. Θεωρήσαμε πως ήταν καλύτερα να αποφευχθούν επιπλέον εξαρτήσεις ειδικά όταν είναι η cool database du jour.

# Date parser

Η κατανόηση ημεφομηνιών υλοποιήθηκε σε ένα ξεχωφιστό πακέτο που ονομάζεται overlay-parse[22].

#### 1. Parsing με overlays

Η έννοια του overlay εμπνεύστηκε από τα emacs overlays[6]. Είναι αντικείμενα που εξειδικεύουν την συμπεριφορά ενός υποσυνόλου του κειμένου με το να του δίνουν ιδιότητες για παράδειγμα το κάνουν clickable ή highlighted.

Ένα overlay επί ενός μέρους ενός κείμενου t στο πλαίσιο μας είναι:

- Ένα ζευγάρι φυσικών αριθμών που ορίζει την έκταση του υποκείμενου
- ένα σύνολο από ετικέτες (tag set) που ορίζουν τα εννοιολογικά σύνολα στα οποία εμπίπτει το συγκεκριμένο υποκείμενο.
- Αυθαίφετες πληφοφοφίες (τύπου) που το συγκεκφιμένο υποκείμενο εκφράζει.

Πιο αυστηρά:

$$o_i \in TextRanget \times Set(Tag) \times A$$
 numbers  
 $Text \rightarrow \{o_1, o_2, ..., o_n\}$ 

Για παράδειγμα, από το παρακάτω κείμενο

The weather today, 
$$\overbrace{Tuesday}^{o_1}$$
  $\overbrace{21}^{st}$  of  $\overbrace{November}^{o_3}$   $\overbrace{2016}^{o_4}$ , was sunny.

Μπορούμε να εξάγουμε overlays  $\{o_1,...,o_4\}$  έτσι ώστε

$$o_1 = (r"Tuesday"),$$
 {DayOfWeek, FullName}, 2)  
 $o_2 = (r"21^{st"}),$  {DayOfMonth, Numeric}, 21)  
 $o_3 = (r"November"),$  {Month, FullName}, 11)  
 $o_4 = (r"2016"),$  {Year, 4digit}, 2016)

Παρατηρούμε ότι όλα τα overlays του παραδείγματος έχουν  $A=\mathbb{N}$ , όπως κωδικοποιούμε την ημέρα της εβδομάδος, τη μέρα του μήνα, το μήνα του έτους ως φυσικούς αριθμούς. Κωδικοποιούμε πιο ακριβή πληροφορία (πχ αυτή η μέρα είναι διαφορετική από το μήνα από την φύση της) στο σύνολο των ετικετών (tag sets).

Όταν έχουμε ένα σύνολο από overlays μπορούμε να ορίσουμε overlay sequences ως overlays τα οποία είναι κατά συνέχεια, Αυτά και τα δικά τους tag sets ταυτίζονται με ειδικά μοτίβα. Για παράδιγμα μπορούμε να ψάξουμε για σειρές από overlays που ταιριάζουν με το pattern

$$p = \text{DayOfMonth}, \text{Separator}(/), (\text{Month} \land \text{Number}), \text{Separator}(/), \text{Year}$$

ταιοιάζει patterns όπως 22/07/1991, οπού Separator(/) ταιοιάζει μονό με τον χαραχτήρα "/"

#### 2. Το παράδειγμα των ημερομηνιών

Η βασική εφαρμογή που θα χρησιμοποιήσουμε ως παράδειγμα για τη λειτουργία των overlays είναι η κατανόηση ημερομηνιών. Το dates sumbmodule έχει 2 βασικά entry points:

- just\_dates που ψάχνει για ημεφομηνίες σε ένα κείμενο.
- just\_ranges που ψάχνει για εύρη ημερομηνιών σε ένα κείμενο.

Παρακάτω παρουσιάζονται κάποια παραδείγματα. Σημειώστε πως 0 σημαίνει unspecified

```
>>> from overlay_parse.dates import just_dates,
    just_ranges, just_props
>>> just_dates("Timestamp:_{\square}22071991:_{\square}She_{\square}said_{\square}she_{\square}
\verb|uuuuuuu| cominguonuaprilutheu18th, uit'su26uapru2014u
    and \( \text{hope} \( \text{is} \) leaving \( \text{me."} \)
\dots [(22, 7, 1991), (18, 4, 0), (26, 4, 2014)]
>>> dates = just_dates("200 \sqcup AD \sqcup 300 \sqcup b.c.")
>>> just_dates("200_\AD_\300_\b.c.")
[(0, 0, 200), (0, 0, -300)]
>>> just_ranges(u"I_will_be_there_from_2008_to_2009"
[((0, 0, 2008), (0, 0, 2009))]
>>> just_ranges("I_{\sqcup}will_\sqcupstay_{\sqcup}from_{\sqcup}July_{\sqcup}the_{\sqcup}20th_{\sqcup}
    until utoday")
[((20, 7, 0), (29, 4, 2016))]
>>> just_dates('{{Birthudateuanduage|1969|7|10|df=y
    }}')
[(10, 7, 1969)]
>>> just_ranges(u'German: | u02c8vu0254lfu0261a |
    \verb"u014b|| \verb"ama|| \verb"u02c8de|| \verb"u02d0|| \verb"u028as|| \verb||| \verb"u02c8mo|| \verb"u02d0tsa||
    \u0281t], \u0281t], \u0281t], \u0281t], \u0281t], \u0281t], \u0281t], \u0281t]
    xa0\u2013_{\sqcup}5_{\sqcup}December_{\sqcup}1791'
[((27, 1, 1756), (5, 12, 1791))]
```

# Παραρτήματα

1. Παράδειγμα python unit test

```
class TestFetcher(unittest.TestCase):
    def setUp(self):
        self.fetcher = fetcher.get_fetcher()
   def test_html(self):
        html = self.fetcher.html_source("Led_
           Zeppelin")
        self.assertIn("Jimmy⊔Page", html)
    def test_markup_source(self):
        src = self.fetcher.markup_source("Led_
           Zeppelin")
        self.assertIn("{{Infobox⊔musical⊔artist",
   def test_unicode_html(self):
       html = self.fetcher.html_source(u"Rhône")
        self.assertIn("France", html)
    def test_unicode_source(self):
        src = self.fetcher.markup_source("Rhône")
        self.assertIn("Geobox|River", src)
    def test_silent_redirect(self):
        # redirects are only supported when
           force_live is set to True
        src = self.fetcher.markup_source("Obama",
```

force\_live=True)
self.assertFalse(re.match(fetcher.
 REDIRECT\_REGEX, src))

#### 2. Παράδειγμα εκτέλεσης ενός python test

```
$ python setup.py test -s tests.test_lispify
running test
running egg_info
writing requirements to wikipediabase.egg-info/
   requires.txt
writing wikipediabase.egg-info/PKG-INFO
writing top-level names to wikipediabase.egg-info/
   top_level.txt
writing dependency_links to wikipediabase.egg-info/
   dependency_links.txt
writing entry points to wikipediabase.egg-info/
   entry_points.txt
reading manifest file 'wikipediabase.egg-info/
   SOURCES.txt'
reading manifest template 'MANIFEST.in'
writing manifest file 'wikipediabase.egg-info/
   SOURCES.txt'
running build_ext
test_bool (tests.test_lispify.TestLispify) ... ok
test_bool_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_date_multiple_voting (tests.test_lispify.
   TestLispify) ... ok
test_date_simple (tests.test_lispify.TestLispify)
   ... ok
test_date_with_range (tests.test_lispify.TestLispify
   ) ... ok
test_dict (tests.test_lispify.TestLispify) ... ok
test_dict_with_escaped_string (tests.test_lispify.
   TestLispify) ... ok
test_dict_with_list (tests.test_lispify.TestLispify)
    ... ok
test_double_nested_list (tests.test_lispify.
   TestLispify) ... ok
test_error (tests.test_lispify.TestLispify) ... ok
test_error_from_exception (tests.test_lispify.
   TestLispify) ... ok
test_keyword (tests.test_lispify.TestLispify) ... ok
```

```
test_keyword_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_list (tests.test_lispify.TestLispify) ... ok
test_list_of_dict (tests.test_lispify.TestLispify)
test_list_of_dict_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_list_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_nested_list (tests.test_lispify.TestLispify)
   ... ok
test_none (tests.test_lispify.TestLispify) ... ok
test_none_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_number (tests.test_lispify.TestLispify) ... ok
test_number_with_typecode (tests.test_lispify.
   TestLispify) ... ok
{\tt test\_string\ (tests.test\_lispify.TestLispify)\ \dots\ ok}
test_string_escaped (tests.test_lispify.TestLispify)
    ... ok
test_string_not_keyword (tests.test_lispify.
   TestLispify) ... ok
test_string_with_typecode (tests.test_lispify.
   TestLispify) ... ok
test_unicode_string (tests.test_lispify.TestLispify)
Ran 27 tests in 0.047s
OK
```

#### 3. Βοίσκοντας συνώνυμα με MySQL

```
mysql> select page_title, rd_title from \
redirect join page on
rd_from = page_id and
(rd_title = "Bill_Clinton" or page_title = "
    Bill_Clinton");
(rd_title = "Bill_Clinton" or page_title = "
    Bill_Clinton");
+--
```

#### 46 rows in set (11.77 sec)

_page_title	rd_title
BillClinton	Bill_Clinton
William_Jefferson_Clinton	Bill_Clinton
President_Clinton	Bill_Clinton
William_Jefferson_Blythe_IV	Bill_Clinton
Bill_Blythe_IV	Bill_Clinton
Clinton_Gore_Administration	Bill_Clinton
<pre>Buddy_(Clinton's_dog)</pre>	Bill_Clinton
Bill_clinton	Bill_Clinton
William_Jefferson_Blythe_III	Bill_Clinton
President_Bill_Clinton	Bill_Clinton
Bull_Clinton	Bill_Clinton
Clinton,_Bill	Bill_Clinton
William_clinton	Bill_Clinton
42nd_President_of_the_United_States	Bill_Clinton
Bill_Jefferson_Clinton	Bill_Clinton
William_JClinton	Bill_Clinton
Billl_Clinton	Bill_Clinton
Bill_Clinton\	Bill_Clinton
Bill_Clinton's_Post_Presidency	Bill_Clinton
Bill_Clinton's_Post-Presidency	Bill_Clinton

Continued on next page

#### Continued from previous page

- Continued from provious page	
page_title	rd_title
Klin-ton	$Bill_Clinton$
Bill_JClinton	Bill_Clinton
William_Jefferson_"Bill"_Clinton	Bill_Clinton
${\tt William\_Blythe\_III}$	Bill_Clinton
${\tt William\_J.\_Blythe}$	Bill_Clinton
${\tt William\_J.\_Blythe\_III}$	Bill_Clinton
Bil_Clinton	Bill_Clinton
WilliamJeffersonClinton	$Bill_Clinton$
William_J_Clinton	Bill_Clinton
Bill_Clinton's_sex_scandals	Bill_Clinton
Billy_Clinton	Bill_Clinton
Willam_Jefferson_Blythe_III	$Bill_Clinton$
William_"Bill"_Clinton	Bill_Clinton
Billll_Clinton	Bill_Clinton
Bill_Klinton	Bill_Clinton
William_Clinton	Bill_Clinton
Willy_Clinton	Bill_Clinton
${\tt William\_Jefferson\_(Bill)\_Clinton}$	$Bill_Clinton$
Bubba_Clinton	Bill_Clinton
MTV_president	Bill_Clinton
MTV_President	Bill_Clinton
The_MTV_President	Bill_Clinton
Howard_GPaster	Bill_Clinton
Clintonesque	Bill_Clinton
William_Clinton	Bill_Clinton
William_Jefferson_Clinton	Bill_Clinton

#### 4. Χαραμτηριστικά για τον person classifier

#### (a) Category regexes

Χοησιμοποιεί τις απόλουθες συνήθεις επφράσεις (regular expressions) για να ταυτίσει τις πατηγορίες ενός άρθρου.

- .\* person
- \^\d+ deaths.\*
- \^\d+ births.\*
- .\* actors
- .\* deities

- .\* gods
- .\* goddesses
- .\* musicians
- .\* players
- .\* singers

#### (b) Category regex excludes

Αποκλείει τις ακόλουθες regexes.

- \sbased on\s
- \sabout\s
- lists of\s
- animal\s

#### (c) Category matches

Γνωρίζουμε ότι ένα άρθρο αναφέρεται σε ένα πρόσωπο εάν η σελίδα ανήκει σε μια από τις ακόλουθες mediawikia κατηγορίες:

- american actors
- american television actor stubs
- american television actors
- architects
- british mps
- · character actors
- computer scientist
- dead people rumoured to be living
- deities
- disappeared people
- fictional characters
- film actors
- living people
- musician stubs
- singer stubs
- star stubs
- united kingdom writer stubs
- united states singer stubs
- writer stubs
- year of birth missing

• year of death missing

Για ένα παράδειγμα δείτε το παράρτημα.

Όπως είναι φανεφό η λίστα με τις κατηγοφίες είναι αθυαίφετη και όχι πλήφης. Πολλαπλές μέθοδοι μποφούν να χρησιμοποιηθούν για να διοφθωθεί αυτό. Μεφικές από αυτές είναι:

- Μέθοδοι με Supervised machine learning όπως SVM χρησιμοποιώντας άλλες μεθόδους να ορίσουν ένα πρόσωπο και να δημιουργήσουν εκπαιδευτικές ομάδες.
- Εμπλουτίζοντας την υπάρχουσα λίστα κατηγοριών χρησιμοποιώντας στατιστικά από κατηγορίες άρθρων που έχουμε βρει με άλλους τρόπους ότι αναφέρονται σε πρόσωπα.

#### 5. Παράδειγμα κατηγοριών άρθρων

Το άρθοο που αναφέρεται στον Leonardo DiCaprio εντασσεται στις επόμενες κατηγορίες (με bold είναι η κατηγορία που χρησιμοποιεί το Wikipedia-Base για να αποφασίσει πως το άρθοο αναφέρεται σε άνθρωπο).

- Leonardo DiCaprio
- 1974 births
- Living people
- 20th-century American male actors
- 21st-century American male actors
- American environmentalists
- American film producers
- · American male child actors
- American male film actors
- American male soap opera actors
- American male television actors
- American people of German descent
- American people of Italian descent
- American people of Russian descent
- American philanthropists
- Best Actor AACTA Award winners
- Best Actor Academy Award winners

- Best Drama Actor Golden Globe (film) winners
- Best Musical or Comedy Actor Golden Globe (film) winners
- California Democrats
- Film producers from California
- Formula E team owners
- Male actors from Hollywood, California
- Male actors from Palm Springs, California
- Male actors of Italian descent
- People from Echo Park, Los Angeles
- Silver Bear for Best Actor winners

Οι κατηγορίες αυτές μοιάζουν ως εξής στο άρθρο

Categories: Leonardo DiCaprio | 1974 births | Living people | 20th-century American male actors | 21st-century American male actors | American enterior | American male actors | Americ

Figure 9.1: The rendered list of categores for Leonardo DiCaprio

# Part III WikipediaMirror

Wikipedia mirror είναι ένα σύστημα με στόχο να αυτοματοποιήσει τη δημιουργία ενός τοπικού κλώνου της wikipedia περιέχοντας μόνο τα άρθρα — δεν περιέχει τους χρήστε , συζήτηση και ιστορικό επεξεργασιών. Η αυτοματοποιημένη διαδικασία περιλαμβάνει τη ρύθμιση ενός διακομιστή, μια βάση δεδομένων και γέμισμα αυτής της βάσης δεδομένων με τα άρθρα της wikipedia Ο σκοπός της είναι να παρέχει την δυνατότητα πρόσβασης του συνόλου των δεδομένων της Wikipedia, ανεξάρτητα από το wikipedia.org.

# Mediawiki stack overview

Το wikipedia-mirror βασίζεται στο MediaWiki stack που παφέχεται από το Bitnami, μια υπηφεσία που χτίζει το σύνολο του διακομιστή εντός των οφίων ενός direcotry. Αυτό είναι χφήσιμο γιατί αποφεύγεται η επιβάφυνση μέσω της χφήσης ενός container ή VM τεχνολογίας και μας δίνει τη δυνατότητα να έχουμε άμεση πφόσβαση στο σύστημα αφχείων του stack, ενώ εξακολουθούμε να έχουμε το σύστημα κατασκευής Bitnami να κάνει την κοπιώδη εφγασία της ενοφχήστφωσης των διαφόφων τμημάτων και επίσης διαχωφίζεται ο διακομιστής από το υπόλοιπο του συστήματος.

Το stack αποτελείται από

- Έναν http server, στην περίπτωση μας τον apache [3]
- Ένα web application runtime, στην περίπτωση μας PHP[23]
- Μια βάση δεδομένων, στην περίπτωση μας η MySQL
- Το ίδιο το web application, δηλαδή mediawiki

Όλα τα παραπάνω παρέχονται από το bitnami mediawiki stack. Το Xampp[33] παλιότερα ήταν αποδεκτά η καλύτερη επιλογή αλλά είναι unmaintained, έτσι αποφασίσαμε να χρησιμοποιήσουμε το bitnami το οποίο δουλεύει αρκετά καλά.

Όταν το stack ουθμιστεί κατάλληλα, το wikipedia dump xml κατεβαίνει και μετατρέπεται σε sql dump με mwdumper[17]. Θα μπορούσε να κάνουμε pipe στο MySQL αλλά η εξαγωγή είναι χρονοβόρα και είναι πιθανό να προκύψουν προβλήματα κατά το dumping.

1. Στοιχεία του stack.

Παρουσιάζεται κάθε στοιχείο του stack με περισσότερες λεπτομέρειες παρακάτω.

#### (a) Apache

Σύμφωνα με τη wikipedia:

The Apache HTTP Server, colloquially called Apache, is the world's most used web server software. Originally based on the NCSA HTTPd server, development of Apache began in early 1995 after work on the NCSA code stalled. Apache played a key role in the initial growth of the World Wide Web, quickly overtaking NCSA HTTPd as the dominant HTTP server, and has remained most popular since April 1996. In 2009, it became the first web server software to serve more than 100 million websites.

Apache is developed and maintained by an open community of developers under the auspices of the Apache Software Foundation. Most commonly used on a Unix-like system (usually Linux), the software is available for a wide variety of operating systems besides Unix, including eComStation, Microsoft Windows, NetWare, OpenVMS, OS/2, and TPF. Released under the Apache License, Apache is free and open-source software.

Είναι δίκαιο να πούμε ότι ο apache είναι ένας από τους πιο δημοφιλείς διακομιστές web στο διαδίκτυο. Η ίδια η wikipedia.org φαίνεται να χρησιμοποιεί ένα πιο σύνθετο stack που περιλαμβάνει τον varnish, έναν HTTP επιταχυντή, και nginx[20], μια εναλλακτική του apache, επίσης αρκετά δημοφιλή διακομιστή HTTP. Καταλήξαμε σε αυτό το συμπέρασμα από την επιθεώρηση των headers που επιστρέφονται από τη wikipedia.org. Στην περίπτωση http://www.wikipedia.org ανακατευθυνόμαστε προς το secure domain (προσοχή στη γραμμή Server:):

```
$ curl -s -D - http://www.wikipedia.org -o /dev/
    null
HTTP/1.1 301 TLS Redirect
Server: Varnish
[...]
```

Και αν ζητήσουμε κατ ευθείαν για το https://www.wikipedia.org

```
$ curl -s -D - https://www.wikipedia.org -o /dev
    /null
HTTP/1.1 200 OK
Server: nginx/1.9.4
[...]
```

Ωστόσο, είναι πέρα από το πεδίο της συγκεκριμένης εργασίας να αναπαράγουμε με ακρίβεια την υποδομή της Wikipedia. Έχουμε αποκλειστικά επικεντρωθεί στην λειτουργικότητά της. Γι αυτό, λόγω της δημοτικότητας, και της ταχύτητας της αυτόματης εγκατάστασης του Bitnami MediaWiki stack χρησιμοποιήθηκε ως διακομιστή μας ο apache.

#### (b) PHP

Η MediaWiki , η οποία συζητείται παρακάτω, είναι γραμμένη εξ ολοκλήρου σε PHP, μια δημοφιλή server-side γλώσσα προγραμματισμού, με dynamic typing, object-oriented scripting γλώσσα. Η PHP εγκαθίσταται με το Bitnami mediawiki stack. Η PHP είναι δημοφιλής ανάμεσα στους προγραμματιστές του web και αυτό οφείλεται εν μέρει στην υποστήριξη που έχει από πολλές σχετικές βιβλιοθήκες για βάσεις δεδομένων (συμπεριλαμβανομένων PostgreSQL, MySQL Microsoft SQL Server και SQLite) και είναι ουσιαστικά ένα template δημιουργίας προτύπων γλώσσας HTML.

#### (c) MySQL

Mediawiki μποφεί να χφησιμοποιήσει πληθώφα SQL database backends:

- MSSQL: Μια SQL βάση από τη Microsoft
- MySQL: Χοησιμοποιώντας τη standard PHP library για MySQL.
- MySQLi: Μια επέκταση του MySQL backend
- Oracle: Μια αποκλειστικής εκμεταλλεύσεως SQL database από την Oracle.
- SQLite: Μια SQL database που συνήθως χοησιμοποιείται ως βιβλιοθήκη αντί για client-server scheme όπως γίνεται με τις άλλες επιλογές της λίστας

Η Wikipedia παφέχει πολλαπλά dump files για τους SQL πίνακες δευτεφεύοντος σημασίας στο MySQL format (eg. page redirects, categories etc) και προτείνει mwdumper που μετατφέπει τα XML dumpls των άφθρων της wikipedia σε MySQL. Το γεγονός αυτό,

και το ότι παρέχεται με το αυτοματοποιημένο stack του bitnami, κάνει τη MySQL την προφανή επιλογή για το wikipedia-mirror stack.

#### (d) Mediawiki

To Mediawiki είναι η μαρδιά της wikipedia. Είναι ένα free μαι open-source[7] wiki application. Δημιουργήθηκε από το Wikimedia Foundation[28] μαι τρέχει πολλά δημοφιλή site όπως Wikipedia, Wikitionary[32] μαι Wikimedia Commons[29].

Το λογισμικό έχει περισσότερες από 800 ρυθμίσεις και περισσότερες από 2.000 επεκτάσεις διαθέσιμες για τη διευκόλυνση προσθήκης ή αλλαγής διάφορων χαρακτηριστικών. Αποκλειστικά στη Wikipedia, πάνω από 1000 αυτοματοποιημένα και ήμι-αυτοματοποιημένα bots και άλλα εργαλεία έχουν αναπτυχθεί για να βοηθήσουν στο moderation. Τα περισσότερα από αυτά δεν έχουν σημασία για τους δικούς μας σκοπούς. Οι χρήσιμες για μας επεκτάσεις είναι οι Scriunto και parserfunctions, και οι μόνες χρήσιμες ρυθμίσεις σχετίζονται με το όνομα της τοποθεσίας, το όνομα της βάσης δεδομένων κλπ. Ως επί το πλείστον αυτές τις διαχειρίζεται το Bitnami.

# **Setting up**

Στη συνέχεια είναι βήμα προς βήμα οδηγείες για να στήσει κάνεις το wikipedia mirror. Πρώτα κατεβάζουμε τον κωδικά χρησιμοποιοντας το git[8]:

```
$ git clone https://github.com/fakedrake/wikipedia-
    mirror
$ cd wikipedia-mirror
```

Σ' αυτό το σημείο θεωρητικά κάποιος μπορεί να τρέξει make sql-load-dumps τα οποία θα φροντίσουν να στηθεί οτιδήποτε χρειάζεται σε μορφή dumps για να φορτωθεί σε μια λειτουργική SQL βάση δεδομένων. Φυσικά για να γίνει αυτό πρώτα θα εκτελεσθούν μερικά βήματα.

- Να κατεβάσουμε τα wikipedia database dumps σε XML format.
- Να τα μετατρέψουμε σε format που καταλαβαίνει η MySQL.
- Να στήσουμε το bitnami stack που περιλαμβάνει ένα local install της MySQL
- Να φορτώσουμε τα MySQL dumps στη MySQL.

Όλα αυτά τα βήματα κωδικοποιούνται ως τμήμα μιας ιεραρχίας εξαρτήσεων σε makefile targets και στη θεωρία αυτό πραγματοποιείται αυτόματα και δημιουργεί μια wikipedia mirror. Όμως αυτή λειτουργία είναι μεγάλη και εύθραυστη και συνιστάται κάθε βήμα να γίνεται εξατομικευμένα και χειροκίνητα.

Ποώτα, κατεβάζουμε και εγκαθιστάμε το bitnami. Η ακόλουθη εντολή θα κατεβάσει ένα executable από το bitnami website και θα κάνει μια τοπική εγκατάσταση του bitnami stack όπως συζητήθηκε παραπάνω:

#### \$ make bmw-install

Το επόμενο βήμα είναι να βεβαιωθούμε ότι το maven, ένα εργαλείο για software project management για java, είναι εγκαταστημένο. Απαιτείται για να εγκατασταθεί και να στηθεί το mwdumper (βλέπε παρακάτω). Μπορεί να γίνει αυτό αν βεβαιωθούμε ότι το παρακάτω δουλεύει:

\$ mvn --version

Σημείωση: Αν είμαστε σε Ubuntu 14.04, μπορούμε να εγκαταστήσουμε το Maven (για Java) χρησιμοποιώντας sudo apt-get install maven.

Τώρα όλα είναι έτοιμα για το αυτόματο download Wikipedia's XML dumps[30] και στη συνέχεια να μετατρέπει σε SQL χρησιμοποιώντας mwdumper. Πρώτα το mwdumper θα πρέπει να κατέβει και να χτισθεί. Μετά τα συμπιεσμένα XML dumps θα πρέπει να κατέβουν από την wikipedia. Θα αποσυμπιεστούν και τελικά θα μετατραπούν σε MySQL dumps χρησιμοποιώντας mwdumper. Αυτή είναι πολύ χρονοβόρα διαδικασία και χρειάζεται 6-11 ώρες σε ένα τυπικό μηχάνημα:

\$ make sql-dump-parts

Όταν αυτο γίνει επιτυχώς μπο<br/>ρούμε να φορτώσουμε τα SQL dumps στη βάση δεδομένων MySQL

\$ make sql-load-parts

Και τελικά

\$ make mw-extensions

Για να ουθμιστούν τα mediawiki extensions.

### Mediawiki extensions

Για να ενεργήσει η MediaWiki όπως η wikipedia απαιτούνται μια σειρά από extensions. Η διαδικασία εγκατάστασης των εν λόγω extensions δεν είναι αυτοματοποιημένη. Για να γίνει αυτόματη διαχείριση αυτής της πολύπλοκης διαδικασίας παρέχεται ένας μηχανισμός για την εγκατάσταση των extensions. Για να υποστηρίξουμε επιπλέον extensions για την wikipediabase πρέπει να προσθέσουμε τον ακόλουθο κώδικα στο Makefile.mwextnesions (τροποποιημένο αναλόγως)

```
MW_EXTENSIONS += newextension
mw-newextension-url = url/to/new/extnesion/package.tar.
    gz
mw-newextension-php = NewExtensionFile.php
mw-newextension-config = '$$phpConfigVariable = "value"
    ;'
```

Η wikipedia-mirror θα φουντίσει ώστε το extension να είναι ήδη εγματαστημένο και εάν δεν είναι θα τοποθετήσει τα σωστά αρχεία στο σωστό μέρος και θα διορθώσει τους ματάλληλους configuration files. Τα entry points για την διαχείριση των extensions είναι (με την προϋπόθεσή ότι το όνομά του εγγραφομένων extensions είναι newextension):

```
make mw-print-registered-extensions # Output a list of the registed extensions
make mw-newextension-enable # Install and/or enable the extension
make mw-newextension-reinstall # Reinstall an extension
```

make mw-newextension-disable extensionmake mw-newextension-clean # Remove the extension

# Disable the

Όλα τα registered extensions θα εγκατασταθούν και θα ενεργοποιηθούν όταν το wikipedia-mirror έχει χτισθεί.

# Φορτώνοντας τα mediawiki dumps

Η Wikipedia παρέχει μηνιαία dumps όλων των βάσεων δεδομένων της. Το μεγαλύτερο μέρος των dumps είναι σε μορφή ΧΜL και πρέπει να κωδικοποιούνται σε MySQL για να φορτωθούν στη βάση δεδομένων του wikipedia-mirror. Υπάρχουν περισσότεροι από ένας τρόποι να το κάνουμε αυτό.

Το Mediawiki παπετάρεται με ένα βοηθητιπό πρόγραμμα για την εισαγωγή του XML dump. Ωστόσο, η χρήση του για την εισαγωγή ενός πλήρους wikipediamirror αποθαρρύνεται λόγω των περιορισμών των επιδόσεων. Αντ αυτού προτείνονται εργαλεία όπως το mwdumper που μετατρέπουν τα XML dumps σε MySQL ερωτήματα τα οποία γεμίζουν τη βάση δεδομένων.

Το mwdumper είναι γραμμένο σε Java και αποστέλλεται χωριστά από το MediaWiki και μπορεί να μετατρέψει τα δεδομένα μεταξύ των ακόλουθων μορφών:

- XML
- MySQL dump
- SQLite dump
- CSV

Για τους σκοπούς μας έχει ενδιαφέρον μονό ο μετασχηματισμός από XML σε MySQL, ωστόσο συναντήθηκαν σημαντικές δυσκολίες σε αυτή τη διαδικασία. Λεπτομέρειες για το ποιές ήταν και πως αντιμετοπίστηκαν δείτε την περιγραφή του xerces bug στο παράρτημα.

# Εργαλεία

Ένας αφιθμός εφγαλείων αναπτύχθηκε για να βοηθήσουν τη διαδικασία του χειφισμού και της παφακολούθησης της διαδικασία του φοφτώματος των dumps στη βάση δεδομένων. Παφουσιάζονται με λεπτομέφεια παφακάτω. Εφ όσον ο πηγαίος κώδικάς τους είναι συνοπτικός παφατίθεται ολόκληφος στο παφάφτημα

#### 1. utf8thread.c

To utf8thread.c είναι ένα χαμηλού επιπέδου πρόγραμμα το οποίο γεμίζει με κενά όλα τα invalid utf-8 characters από το αρχείο. Χρησιμοποιούμε pthreads για να επιταχύνουμε τα πράγματα.

#### 2. webmonitor.py

Το webmonitor. py είναι ένα python script το οποίο σερβίσει μια σελίδα που δείχνει live δεδομένα σε μοφή ιστογράμματος για την πρόοδο του γεμίσματος της βάσης δεδομένων. Το webmonitor. py σερβίσει στατικές html σελίδες και μετά τους στέλνει δεδομένα μέσω websocket. Το Webmonitor μπορεί να δείχνει οποιοδήποτε stream από ζευγάρια <epoc date> <float value> που λαμβάνει στο input. Για παράδειγμα:

\$ pip install tornado

Ποώτα ποέπει να εγματαστήσουμε τα dependencies του script. Το οποίο μπορεί να είναι tornado, anasynchronous web framework που υποστηρίζει websockets. Δίνουμε οδηγίες στο tornado, που θα σερβίρει την ακόλουθη σελίδα:

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD_HTML_4.01//EN" "</pre>
   http://www.w3.org/TR/html4/strict.dtd">
<html>
  <head>
    <meta http-equiv="Content-Type" content="text/</pre>
       html; _charset=utf-8">
    <title>DrNinjaBatmans Websockets</title>
    <script type="text/javascript" src="http://code.</pre>
        jquery.com/jquery-1.10.1.js"></script>
    <script type="text/javascript" src="http://code.</pre>
       highcharts.com/highcharts.js"></script>
    <script>
     var chart; // global
     var url = location.hostname + ':' + (parseInt(
        location.port));
     var ws = new WebSocket('ws://' + url + '/
        websocket');
     ws.onmessage = function(msg) {
         add_point(msg.data);
     };
     // ws.onclose = function() { alert('Connection
        closed.'); };
     var add_point = function(point) {
         var series = chart.series[0],
         shift = series.data.length > %d;
         chart.series[0].addPoint(eval(point), true,
              shift);
     };
     $(document).ready(function() {
         chart = new Highcharts.Chart(JSON.parse('%s
             '));
     });
    </script>
  </head>
  <body>
      <div id="container" style="width:\(\_800px;\)\(\)</pre>
         height: \_400px; \_margin: \_0\_auto"></div>
```

```
</body>
```

Η σελίδα αναμένεται να διαβάζει ένα stream τιμών από ένα websocket στο ws://localhost:8888/hostname και τα κάνει plot σε πραγματικό χρόνο χρησιμοποιώντας το highcharts.js. Ο προσεκτικός αναγνώστης ίσως παρατηρήσει ότι τα παραπάνω δεν είναι ακριβώς HTML αλλά περισσότερο ένα python formatted string. Αυτό συμβαίνει για 2 λόγους. Πρώτον γιατί το python script διαχειρίζεται το configuration και δεύτερον γιατί το πλάτος του graph υπολογίζεται σε page load time και το plot χρειάζεται να μετατοπισθεί για να δείξει μόνο τα πιο πρόσφατα σημεία.

```
$ for i in \{1..100\}; do echo $i; sleep 1; done | \ awk -oL "\{print_{\square}\1/100}" | \ python webmonitor.py
```

Αυτό θα παράγει σε διάστημα 1 δευτερολέπτου, αριθμούς από το 1 ως το 100. Μετά θα τα κανονικοποιήσει χρησιμοποιώντας awk και τα θα τα τροφοδοτήσει στο webmonitor. Αφού αυτή η εντολή εκτελεσθεί μπορούμε να ανοίξουμε τον browser και να κάνουμε navigate στο localhost:8888.

Χοησιμοποιούμε αυτό το script για να ελέγξουμε από απόσταση το ολικό μέγεθος των δεδομένων που το mysql λαμβάνει, μιας που αυτή είναι η πιο χρονοβόρα διαδικασία της δημιουργίας του mirror.

#### 3. xml-parse.sh

Άπλα αφαιρώντας συγκεκριμένα άρθρα διορθώνουμε το πρόβλημα του xerces. Αν τα άρθρα είναι απομονωμένα τότε επίσης το error εξαφανίζεται. Το xml-parse.sh διαγράφει τα ζητούμενα άρθρα από το αρχείο xml.

```
xml-parse.sh <original-xml-file> <
    title_of_article_to_remove> [inplace]
```

αν το τελευταίο όρισμα είναι το inplace, τότε το page\_remover.c θα χρησιμοποιηθεί για να καλύψει το άρθρο με κενά. Αυτή η διαδικασία είναι πολύ πιο γρήγορη. Διαφορετικά η σελίδα άπλα διαγράφεται και το αποτέλεσμα εμφανίζεται στο stdout. Μόλις το xml-parse.sh τελειώσει επιτυχώς μπορεί κανείς να τρέξει:

java -jar tools/mwdumper.jar RESULTING\_XML --format= sql:1.5 > SQL\_DUMP

#### 4. sql-clear.sh

Το sql-clear. sh είναι ένα μικοό bash script που "κουτσουφεύει" όλους τους πίνακες από την βάση δεδομένων. Με τον όφο "κουτσουφευει" εννοούμε ότι αφήνει τα table scheamata ανεπηφέαστα και διαγφάφει όλα τα internal δεδομένα.

#### 5. page\remover.c

Όπως προηγουμένως συζητήθηκε η xerces βιβλιοθήκη την οποία χρησιμοποιεί το mwdumper απέτυχε, φαινομενικά τυχαία, να επεξεργαστεί κάποιες σελίδες. Για να διευθετηθεί αυτό το πρόβλημα αφαιρέσαμε τις σελίδες πλήρως και ξαναπροσπαθήσαμε. Επειδή αυτή η εργασία είναι εύκολη άλλα αργή γράψαμε ένα χαμηλού επιπέδου πρόγραμμα στη C για να το επιλύσουμε, το page\_remover.c. Το page-remover δέχεται ως input το path του XML wikipedia dump, το offset του άρθρου που θέλουμε να καλύψουμε και το μέγεθος του άρθρου. Μετά χρησιμοποιεί το mmap system call για να αποκτήσει ψευδο-random-access στα δεδομένα μέσα στο αρχείο και γεμίζει το άρθρο με withespace characters. Το page\_remover.c δεν είναι threaded μιας που το bottleneck είναι στο HDD IO speed και ο παραλληλισμός δεν θα βοηθούσε.

# Αυτοματισμός

Η δημιουργία μιας wikipedia mirror ίσως φαίνεται μια απλή διαδικασία αλλά συμπεριλαμβάνει πολλές αγκαθωτές λεπτομέρειες και επαναλαμβανόμενα tasks. Πολλαπλές μέθοδοι αυτοματισμού εφαρμόσθηκαν για να ολοκληρώσουν μια μεγάλη ποικιλία tasks που συμπεριλαμβάνονται στην εκτέλεση.

Το πιο σημαντικό μέρος του αυτοματισμού της wikipedia-mirror είναι το make build system. Make είναι ένα build system όπου κάποιος μπορεί να δηλώσει τα απαιτούμενα αρχεία (targets), dependencies για αυτά, και ένα σύνολο από shell commands που θα χτίσουν αυτά τα targets. Κάθε target είναι ουσιαστικά μια finite state machine με δύο καταστάσεις:

- Ένα αρχείο που υπάρχει και είναι επικυροποιημένο με τα dependencies και
- Ένα αρχείο που είτε δεν υπάρχει ή η modification date είναι παλαιότερη από τουλάχιστον ενός εκ των των dependencies.

Επιπλέον πεφιλαμβάνεται μια σειφά από shell εντολές για την μεταφοφά από την πφώτη στη δεύτεφη κατάσταση. Το αποτέλεσμα είναι οτι ο χφήστης απαιτεί τη δημιουφγία ενός αφχείου και το make τφέχει οσο το δυνατόν λιγοτεφες διεφγασίες, αποφεύγοντας να ξαναδημιουφγήσει αφχεία που χφειάζεται αλλα ήδη υπάφχουν. Λεπτομέφειες για τη λειτουφγία των makefiles θα βφείτε στο παφάφτημα.

Χρησιμοποιύμε τα makefiles για να μπορούμε να συνεχίζουμε μια διαδικασία που σταματήσαμε ή απέτυχε οσο πιο κοντα γίνεται στο σημείο οπου εμεινε την τελευταία φορα. Επίσης η λειτουργία των makefiles κανει το make αρκετά έξυπνο ώστε να μην επαναλαμβάνει βήματα που πιθανώς κάναμε χειροκίνητα.

# Επιδόσεις

Το Compile time περιλαμβάνει το χρόνο που χρειάζεται για:

- Κατέβασμα όλων των στοιχείων του wikipedia server
- Το στήσιμο του bitnami stack
  - mwdumper
  - mediawiki-extensions
  - Εγκατάσταση και χτίσιμο αυτών των στοιχείων (~1 min)
  - Κατέβασμα των wikipedia dumps
  - Ποοεπεξεργασία των dumps (~10 mins)
  - Populating τη mysql βάση δεδομένων(~10 days)

Τα Builds έγιναν στο Infolab's Ashmore. Τα system's specs είναι σχετικά ψηλά σε γενικές γραμμές αλλά το bottleneck ήταν το disk IO έτσι λιγότερο από 1% από τις υπόλοιπες διαθέσιμες πηγές χρησιμοποιήθηκαν κατά την διάρκεια του MySQL database population. Συγκεκριμένα τα χαρακτηριστικά του ashmore είναι:

- CPU: Xeon E5-1607 3GHz 4-Core 64 bit
- Main memory: 64G
- HDD: (spinning disk) 500GB + 2Tb

Εφ όσον το βασικό bottleneck είναι η δημιουργία βάσης δεδομένων — δηλαδή οι επιδόσεις της MySQL — δόθηκε μεγάλη προσοχή και πειραματισμός στη σωστή ρύθμιση της βάσης, άλλα η επιτάχυνση ήταν εν τέλη ελάχιστη και έτσι τα περισσότερα απ' όσα δοκιμάστηκαν δεν περιλήφθηκαν στα Makefiles.

Η backend engine που χοησιμοποιήσαμε για τη MySQL είναι η InnoDB. Μερικές από της μεθόδους βελτιστοποίησης που επιχειρήθηκαν παρουσιάζονται παρακάτω.

- Ρύθμιση του innodb\_buffer\_pool\_size[10]. Ενώ η διαθέσιμη μνήμη του ashmore είναι αφκετά μεγάλη, αυξάνοντας το buffer pool μέχρι και κάποια GB δεν είχε σοβαρό αντίκτυπο στην επίδοση.
- Αλλάζοντας το innodb\_flush\_method={{{innodb\_flush\_method}}}
   =0\_DSYNC για να αποφυγυμε κλήσεις στην fsync={{{ref (fsync}}}.
   Ε =fsync είναι ότι ψάχνει σειριακά τις mapped σελίδες ενός αρχείου για dirty pages με αποτέλεσμα να γίνεται αργό για μεγάλα αρχεία.
- Ρυθμίζοντας το innodb\_io\_capacity[12]. Εν τέλη η τιμή της μεταβλητής ήταν υψηλότερη από το bandwidth του σκληρού δίσκου

Η μόνη βελτιστοποίηση που είχε αισθητό αποτέλεσμα ήταν η αλλαγή του MySQL dump ώστε να θέτει

```
SET AUTOCOMMIT = 0; SET FOREIGN_KEY_CHECKS=0;
```

Αυτό επέτρεψε στο InnoDB να κάνει περισσότερη δουλειά στην κύρια μνήμη πριν επικοινωνήσει με το δίσκο και επίσης μείωσε τη συνολική δουλειά εμπιστευόμενη ότι οι τιμές των κελιών που αναφέρονταν σε άλλους πίνακες όντως έδειχναν κάπου.

# Παραρτήματα

#### 1. Το bug στη βιβλιοθήμη xerces

Πιθανότατα η μεγαλύτερη δυσκολία κατά τη δημιουργία του wikipediamirror ήταν η αντιμετώπιση ενός bug στο mwdumper — το εργαλείο για τη μετατροπή των XML dumps σε SQL dumps — το οποίο κάνει το εργαλείο να αποτυγχάνει σε τυχαία άρθρα. Εφ όσον δεν μπορέσαμε να βρούμε τον κριβή λόγο που συμβαίνει αυτό το bug, το παρακάμπτουμε διαγράφοντας τα άρθρα που προκαλούν το πρόβλημα, και εφ όσον είναι ένα μεγάλο εμπόδιο σε μια κατά τα άλλα θεωρητικά πεπατημένη διαδικασία περιγράφουμε τη διαδικασία μας λεπτομερώς.

Ιδού λοιπόν τι ακριβώς συμβαίνει: ενώ τρέχει το make sql-dump-parts εγείρεται το παρακάτω exception:

```
scanContent(Unknown Source)
        at org.apache.xerces.impl.
           {\tt XMLDocumentFragmentScannerImpl\$FragmentContentDispatcher}
           .dispatch(Unknown Source)
        at org.apache.xerces.impl.
           XMLDocumentFragmentScannerImpl.
           scanDocument(Unknown Source)
        at org.apache.xerces.parsers.
           XML11Configuration.parse(Unknown Source)
        at org.apache.xerces.parsers.
           XML11Configuration.parse(Unknown Source)
        at org.apache.xerces.parsers.XMLParser.parse
           (Unknown Source)
        at org.apache.xerces.parsers.
           AbstractSAXParser.parse(Unknown Source)
        at org.apache.xerces.jaxp.
           SAXParserImpl$JAXPSAXParser.parse(
           Unknown Source)
        at javax.xml.parsers.SAXParser.parse(
           SAXParser.java:392)
        at javax.xml.parsers.SAXParser.parse(
           SAXParser.java:195)
        at org.mediawiki.importer.XmlDumpReader.
           readDump(XmlDumpReader.java:88)
        at org.mediawiki.dumper.Dumper.main(Dumper.
           java:142)
make: *** [/scratch/cperivol/wikipedia-mirror/drafts
   /wikipedia-parts/enwiki-20131202-pages-
   articles20.xml-p011125004p013324998.sql] Error 1
```

Διεφευνούμε το πφόβλημα τφέχοντας make --just-print sql-dump-parts με σκοπό να βφούμε επακφιβώς την αλληλουχία εντολών που πφοκάλεσαν το πφόβλημα και ανακαλύπτουμε πως η εντολή που αποτυγχάνει είναι:

```
$ java -jar /scratch/cperivol/wikipedia-mirror/tools
/mwdumper.jar --format=sql:1.5 /scratch/
cperivol/wikipedia-mirror/drafts/wikipedia-parts
/enwiki-20131202-pages-articles20.xml-
p011125004p013324998.fix.xml > /root/path/
wikipedia-parts//enwiki-20131202-pages-
articles20.xml-p011125004p013324998.sql
```

Ευτυχώς αυτή η εντολή δεν τρέχει για μεγάλο διάστημα έτσι μπορούμε

εύχολα να πειραματιστούμε. Εδώ είναι το time output:

```
26.65s user 1.73s system 78% cpu 35.949 total
```

Το λάθος φαίνεται να συμβαίνει κατά την διάρκεια του διαβάσματος του XML dump συνεπώς δεν είναι ειδικό για το SQL output. Αυτό θα μπορούσε να είναι χρήσιμο για να διαπιστώσουμε ποιο άρθρο προκαλεί το λάθος, αποσύροντάς το και ελπίζοντας να λυθεί το πρόβλημα. Για να το εντοπίσουμε κατ αρχάς προσπαθήσαμε να κάνουμε export σε XML:

```
$ java -jar /scratch/cperivol/wikipedia-mirror/tools
/mwdumper.jar --format=xml /scratch/cperivol/
wikipedia-mirror/drafts/wikipedia-parts/enwiki
-20131202-pages-articles20.xml-
p011125004p013324998.fix.xml > /tmp/just-a-copy.
xml
```

Όπως ήταν αναμενόμενο, το ίδιο λάθος εμφανίστηκε. Στη συνέχεια κοιτάμε τα τελευταία δύο άρθρα που έγιναν export. Θέλουμε να μπορούμε να αυτοματοποιήσουμε τη διαδικασία συνεπώς το κάνουμε χρησιμοποιόντας shell commands και όχι με το χέρι: Ξεκινάμε τυπώνοντας με αντίστροφη σειρά το xml αρχείο που δημιουργήθηκε, βρίσκοντας τις τελευταίες δύο εμφανίσεις του <title> χρησιμοποιόντας το εργαλείο grep και αναστρέφοντας ξανά τις προκύπτουσες γραμμές για να τυπώσουμε με την αρχική σειρά. Σημειώνουμε ότι το εργαλείο tac είναι μέρος των coreutils[?]} άλλα όχι του BSD toybox{{{ref(bsd<sub>toybox</sub>}}. Κατά συνέπεια θα υπάρχει σε όλες τις διανομές GNU linux προεγκατεστημένο άλλα όχι και σε mac os x και BSD. Δουλεύει όπως cat, μόνο που εμφανίζει τις γραμμές με αντίστροφη σειρά):

Αυτή η λειτουργία τελειώνει γρήγορα παρά το ότι το /tmp/just-a-copy.xml είναι αρκετά μεγάλο διότι το εργαλείο tac ψάχνει πρώτα το τέλος του αρχείου και διαβάζει προς τα πίσω μέχρι το grep να βρει τα 2

περιστατικά που ψάχνει για <title> και κλείνει. Στο filesystem ext3, και τα παρόμοια filesystems, κατά την διαδικασία αυτή δεν διασχίζεται ολόκληρο το αρχείο. Πράγματι, από τον κώδικα tac βλέπουμε ότι γίνεται χρήση της 1seek που αναζητεί το τέλος του αρχείου χωρίς να το προσπελάσει ολόκληρο και στη συνέχεια διαβάζει αντίστροφα:

```
if (lseek (input_fd, file_pos, SEEK_SET) < 0)</pre>
    error (0, errno, _{-}("%s:\_seek\_failed"), quotef (
       file));
/* Shift the pending record data right to make room
   for the new.
   The source and destination regions probably
      overlap. */
memmove (G_buffer + read_size, G_buffer,
   saved_record_size);
past_end = G_buffer + read_size + saved_record_size;
/* For non-regexp searches, avoid unnecessary
   scanning. */
if (sentinel_length)
    match_start = G_buffer + read_size;
else
    match_start = past_end;
if (safe_read (input_fd, G_buffer, read_size) !=
   read size)
{
    error (0, errno, _("%s:_read_error"), quotef (
       file)):
    return false;
}
```

Ας σώσουμε την διαδοομή από το αρχικό xml αρχείο σε μια μεταβλητή γιατί θα το χρησιμοποιήσουμε χρησιμοποιούμε πολύ. Έτσι από δω και πέρα το \$ORIGINAL\_XML θα έχει διαδρομή από το αρχικό xml.

```
$ export ORIGINAL_XML=/scratch/cperivol/wikipedia-
mirror/drafts/wikipedia-parts/enwiki-20131202-
pages-articles20.xml-p011125004p013324998.fix.
xml
```

Πρώτα ας δούμε αν κάτι περίεργο συμβαίνει με το xml αρχείο

```
\ grep\ "<title>Cranopsis_bocourti</title>" -A 200 -B 100 <math display="inline">\ ORIGINAL\_XML\ |\ less
```

Τίποτα περίεργο δεν βρέθηκε, έτσι δεν μπορούμε πραγματικά να διορθώσουμε το πρόβλημα επιτόπου. Θα προσπαθήσουμε να αποσύρουμε ολόκληρο το άρθρο και ελπίζουμε ότι θα δουλέψει (spoiler alert: δουλεύει). Δεν μπρούμε όμως να κάνουμε κανονικό parsing του xml format μιας που το μέσο μέγεθος του αρχείου που αντιμετωπίζουμε είναι της τάξης των δεκάδες GB. Θα ήταν καλύτερα να βρούμε έναν πιο γρήγορο και πιο low level τρόπο να ανασύρουμε το αρχείο. Θα χρησιμοποιήσουμε λοιπόν καθαρά byte και string operations.

Θα προσπαθήσουμε αρχικά να επιθεωρήσουμε τους parents του τίτλου από το προβληματικό άρθρο. Ευτυχώς το xml που δημιουργήθηκε είναι indented, έτσι μπορούμε να βρούμε τους parents που βασιζόμενοι σε αυτό. Αριθμήσαμε 6 spaces από indentation στη γραμμή που ο mwdumper απέτυχε, έτσι θα ψάξουμε προς τα πίσω από εκεί για τις γραμμές στις οποίες το indentation μειώνεται. Οι γραμμές αυτές θα αντιπροσωπεύουν τους προγόνους του άρθρου μέσα στην ιεραρχία του XML:

```
$ for i in {0..6}; do \
   echo "Level,,$i:"; \
   -m 1 -n | tac; \
done
Level 0:
17564960: <mediawiki xmlns="http://www.mediawiki.org/
   xml/export-0.3/" xmlns:xsi="http://www.w3.org
   /2001/XMLSchema-instance" xsi:schemaLocation="
   http://www.mediawiki.org/xml/export-0.3/uhttp://
   www.mediawiki.org/xml/export-0.3.xsd" version="
   0.3" xml:lang="en">
Level 1:
Level 2:
38: <page>
Level 3:
Level 4:
35:
     <revision>
Level 5:
Level 6:
        <text xml:space="preserve">&lt;!-- This
26:
   article was auto-generated by [[User:Polbot]].
```

Φαίνεται ότι η δομή του συγκεκριμένου xml έχει ως εξής: κάθε page βρίσκεται σε σε ένα domain που ονομάζεται mediawiki. Θα μπορούσαμε να δούμε αυτό επίσης και από την java source αλλά αν και πιο ακριβό από άποψη υπολογιστικών πόρων είναι πιο γρήγορο από το να προσπαθήσουμε να διαβάσουμε τη java του mwdumper.

Ο πιο εύκολος τρόπος να αφαιρέσουμε αυτό το άρθρο θα ήταν το εργαλείο awk. Αλλά είναι πολύ αργό για τους σκοπούς μας και θέλουμε να βελτιστοποιήσουμε και να αυτοματοποιήσουμε την όλη διαδικασία. Πρώτα ας προσπαθήσουμε απλώς να συγκρίνουμε το orginal xml και το xml που δημιουργήσαμε byte προς byte μιας που αυτό είναι πολύ γρήγορη διεργασία, με την ελπίδα πως το άρθρο υπό συζήτηση θα είναι η μόνη διαφορά και έτσι θα μπορούμε πανεύκολα να βρούμε το σημείο οπού θα αρχίσουμε διαγράφουμε:

```
$ cmp /tmp/just-a-copy.xml $ORIGINAL_XML
/tmp/just-a-copy.xml /scratch/cperivol/wikipedia-
mirror/drafts/wikipedia-parts/enwiki-20131202-
pages-articles20.xml-p011125004p013324998.fix.
xml differ: byte 2, line 1
```

Η πρώτη διαφορά είναι στο 20 byte, συνεπώς σίγουρα δεν πρόκειται περί του άρθρου υπό συζήτηση.

```
<namespace key="-1" case="first-letter">
         Special </namespace>
      <namespace key="0" case="first-letter" />
$ head /tmp/just-a-copy.xml
<?xml version="1.0" encoding="utf-8" ?>
<mediawiki xmlns="http://www.mediawiki.org/xml/</pre>
   export-0.3/" xmlns:xsi="http://www.w3.org/2001/
   XMLSchema-instance" xsi:schemaLocation="http://
   www.mediawiki.org/xml/export-0.3/,http://www.
   mediawiki.org/xml/export-0.3.xsd" version="0.3"
   xml:lang="en">
  <siteinfo>
    <sitename>Wikipedia</sitename>
    <base>http://en.wikipedia.org/wiki/Main_Page
    <generator>MediaWiki 1.23wmf4</generator>
    <case>first-letter</case>
    <namespaces>
      <namespace key="-2">Media</namespace>
```

Τα γνωρίσματα των xml tags είναι αρκετά διαφορετικά. Ελπίζουμε τουλάχιστον τα line numbers να συμπίπτουν ώστε αν δε μπορούμε να πάμε κατ ευθείαν στο byte που μας ενδιαφέρει για να αρχίσουμε να διαγράφουμε, τουλάχιστον να μπορούμε γρήγορα να μετρήσουμε τα newlines για να βρούμε το άρθρο. Μετρήσαμε τον αριθμό των γραμμών στο /tmp/just-a-copy.xml και ελπίζουμε ότι η αντίστοιχη γραμμή στο \$ORIGINAL\_XML θα αναφέρεται στο προβληματικό άρθρο. Εάν αυτό συμβεί μπορούμε να αγνοήσουμε τις περιβάλλουσες xml πληροφορίες και να σβήσουμε το προβληματικό άρθρο με βάση αυτήν την πληροφορία. Θα χρησιμοποιήσουμε wc το οποίο είναι αρκετά γρήγορο.

```
$ wc -l /tmp/just-a-copy.xml
17564961 /tmp/just-a-copy.xml
```

Και η αντίστοιχη γραμμή στο \$ORIGINAL\_XML είναι:

```
$ sed "17564960q;d" $ORIGINAL_XML
[[Willie Jones (American football)|Willie Jones]],
```

Ποδόσφαιοο (football)... καμιά σχέση με τα βατράχια (frogs). Φαίνεται ότι δεν μπορούμε να αποφύγουμε κάποιο επίπεδο του parsing.

### (a) Parsing

Θα κάνουμε τις ακόλουθες θεωρήσεις για να αποφύγουμε το σε βάθος parsing του εγγράφου:

- Το XML στο αρχικό αρχείο είναι valid.
- Κάθε ΧΜL μέσα στο άρθρο είναι HTML escaped

Κατ αρχάς δουλεύοντας με γραμμές είναι αργή διαδικασία γιατί το user space code χρειάζεται να ψάχνει newlines. Δουλεύοντας με bytes αναθέτουμε εργασία στο kernel, επιταχύνοντας την εργασία σημαντικά. Έτσι το dd είναι το σωστό εργαλείο για την συγκεκριμένη δουλειά. Αλλά πρώτα θα βρούμε σε ποιο byte είναι το άρθρο που μας ενδιαφέρει

```
$ grep -b "<title>Cranopsis_bocourti</title>" -m
1 $ORIGINAL_XML
1197420547: <title>Cranopsis bocourti</title>
```

Αυτό ίσως πάρει κάποιο χρόνο αλλά δυστυχώς είναι η μόνη μας επιλογή. Η στρατηγική μας είναι να φτιάξουμε 2 αρχεία: το /tmp/original\_tail.xml το οποίο να περιέχει όλα τα δεδομένα που υπάρχουν μετά τη σελίδα που θέλουμε να βγάλουμε και το /tmp/original\_head.xml το οποίο περιέχει όλα τα δεδομένα πριν τη σελίδα που θέλουμε να βγάλουμε.

Τώρα θα χρησιμοποιήσουμε sed να ψάξει για </page> μετά το byte 1197420547 το οποίο θα είναι το σημείο x. Βάζουμε το μέρος του \$ORIGINAL\_XML μετά στο σημείο x μέσα στο αρχείο /tmp/original\_tail.xml:

```
$ dd if=$ORIGINAL_XML skip=1197420547 ibs=1 |
sed '0,/<\/page>/d' > /tmp/original_tail.xml
```

Θαυμάσια, αυτό δούλεψε! Το dd δεν αντιγράφει αντίστροφα έτσι θα χρειαστεί να κάνουμε κάτι πιο περίπλοκο για να κατασκευάσουμε /tmp/original\_head.xml. Ας υποθέσουμε ότι η θέση που βρήκαμε τον τίτλο της σελίδας που θέλουμε να αφαιρέσουμε είναι  $\alpha=1197420547$  και το σημείο που η σελίδα αρχίζει είναι στο σημείο β. Είναι ασφαλές να υποθέσουμε ότι  $\beta>\alpha-1000$  (μπορούσαμε να αλλάξουμε τη σταθερά 1000 εάν αυτή η υπόθεση ήταν λάθος, αλλά τελικά ήταν σωστή). Με αυτό τον τρόπο χρειάζεται μόνο να ψάξουμε στο 1Kb για τη συμβολοσειρά 0Page.

τελικά ενώνουμε το /tmp/original\_head.xml με το /tmp/original\_tail.xml στο αρχείο που δουλεύει κανονικά μ το mwdumper.

## (b) Η τελική λύση

Όλα τα παραπάνω χρησιμοποιήθηκαν για να συντεθεί ένα script που υπάρχει στο data/xml-parse. sh το οποίο χρησιμοποιήθηκε από το makefiles για να απομακρύνει όλα τα προβληματικά άρθρα. Εάν το mwdumper αποτύχει, ταυτοποιούμε το άρθρο που προκάλεσε το πρόβλημα και το απομακρύνουμε χρησιμοποιώντας xml-parse. sh. Στη συνέχεια ξανατρέχουμε το mwdumper. Το επαναλαμβάνουμε αυτό μέχρι το mwdumper να πετύχει. Συνολικά τα προβληματικά άρθρα είναι περίπου 10-15, και είναι διαφορετικά ανάλογα με το dump που χρησιμοποιείται.

#### (c) Καλύπτοντας με κενά

Από την παραπάνω έμθεση των τρόπων που αντιμετωπίσουμε το θέμα του άρθρου που σπάει παραλείψαμε κάτι προφανές. Μια θεματικά διαφορετική προσέγγιση: το να καλύψουμε το άρθρο που προκαλεί το πρόβλημα με κενά. Μόλις εντοπίσουμε το εύρος στο οποίο η σελίδα βρίσκεται μπορούμε να κάνουμε mmap επακριβώς στο τμήμα του \$ORIGINAL\_XML και στη συνέχεια να κάνουμε memset καλύποντας το με χαρακτήρες κενών. Η το πρόγραμμα ζει στο data/page\_remover.c, Παρακάτω παρουσιάζουμε την κλήση στο mmap:

```
ctx->off = off-pa_off;
ctx->fd = open(fname, O_RDWR, 0x0666);
```

```
if (ctx->fd == -1) {
      perror("open");
      return NULL;
 }
 ctx->size = len;
  ctx->data = mmap(0, len+ctx->off, PROT_READ |
     PROT WRITE,
             MAP_SHARED, ctx->fd, pa_off);
 if (ctx->data == MAP FAILED) {
     perror ("mmap");
      return NULL;
 }
και το memset
 /* You MIGHT want to thread this but I dont
     think it will make
   * much more difference than memset. */
 memset(ctx->data + ctx->off, '\( '\)', ctx->size);
```

Περιέργως αυτό δεν έλυσε το πρόβλημα του mwdumper το οποίο δείχνει ότι μάλλον πρόκειται για memory leak από τη μεριά του xerces αλλά αυτό ξεπερνά τους στόχους της παρούσας εργασίας.

## (d) Η εντολή sed

Στο κεφάλαιο για το xerces bug αναφέραμε την χρήση της εντολής sed και ίσως να είναι χρήσιμο να το αναπτύξουμε περεταίρω. sed είναι ένα unix εργαλείο πού βρίσκεται στο πακέτο GNU coreutils το οποίο σύμφωνα με το man page είναι ένας stream editor που φιλτράρει και μεταμορφώνει κείμενο. Η βασική λειτουργία είναι ότι το "pattern space", ή το input stream το οποίο είναι ένα unix stream — που έρχεται από το αρχείο, ένα pipe ή απλά stdin — περνά μέσα από ένα προγραμματίσιμο pipeline. Εκτυπώνεται είτε αυτούσιο το modified pattern space είτε — με τη χρήση του -n flag — επιλεγμένα τμήματα αυτού. Ας δούμε τη χρήση που κάναμε παραπάνω για το sed.

Αρχικά χρησιμοποιήσαμε sed για να εκτυπώσουμε μια μεμονωμένη line σε ένα αρχείο:

```
$ sed "17564960q;d"
```

Αυτό το sed πρόγραμμα διαχωρίζει τις εντολές με semicolon (;). Το sed απαριθμεί τις γραμμές του input stream και τρέχει καθένα από τις εντολές διαχωρισμένες με ";" σε σειρά μέχρι μια να επιτύχει. Οι εντολές εδώ είναι 17564960q και d. Η 17564960q θα σταματήσει το sed όταν φτάσει η παρούσα γραμμή είναι η γραμμή νούμερο 17564960. Η d θα απορρίπτει την παρούσα γραμμή κάθε φορά. Έτσι το sed απορρίπτει γραμμές μέχρι να συναντήσει τη γραμμή 17564960 την οποία εκτυπώνει και τελειώνει.

Χοησιμοποιούμε μια εντολή sed ως μέσος μιας σεισάς εντολών shell piped όλες μαζί με στόχο να εκτυπωθούν όλες ο γραμμές ενός stream μετά από ένα συγκεκριμένο μοτίβο (στην περίπτωσή μας </page>).

αυτή τη φορά είχαμε μόνο μια εντολή sed, d. Το sed απαριθμεί στις γραμμές στο stream, απορρίπτοντας γραμμές στο εύρος των γραμμών από 0 μέχρι τη γραμμή που ταυτίζεται με το <\/page>, ουσιαστικά τυπώνοντας μόνο γραμμές μετά το </page>.

Η τελιχή μας χρήση του sed είναι η αντίστροφη της προηγούμενης,

Εδώ το sed απαφιθμεί ξανά σε όλες τις lines τουstream. Αυτή τη φοφά αποφρίπτοντας γφαμμές ανάμεσα στην πρώτη και αυτήν που ταιφιάζει το <page> μέχρι την τελική γφαμμή, σημειωμένη με \$.

## 2. Λεπτομερώς τα Makefiles

Ας αρχίσουμε με ένα παράδειγμα, σώζουμε το απόλουθο ως Makefile σε ένα project που περιέχει τα αρχεία foo.c, foo.h, bar.c παι bar.h: αυτό σημαίνει ότι για να χτίσουμε το επτελέσιμο foobar χρειαζόμαστε foo.o παι bar.o. Και για να χτίσουμε foo.o παι bar.o χρειαζόμαστε foo.c παι foo.h, παι bar.c παι bar.h αντίστοιχα.

Επίσης παρέχουμε εντολές για να χτισθεί το foo.o, bar.o και foobar, οι οποίες είναι

• gcc foo.c -c -o foo.o

- gcc bar.c -c -o bar.o
- και gcc foo.o bar.o -o foobar

αντίστοιχα. παρατηρούμε ότι δεν υπάρχουν κανόνες για τα .c και .h αρχεία. Αυτό συμβαίνει γιατί το make πρέπει να αποτυγχάνει εάν δεν είναι παρόντα. Έτσι εάν τρέχουμε το make foobar, το make θα ελέγχει για την ύπαρξη του foobar και την ημερομηνία της τροποποίησης. Εάν το foobar λείπει ή η ημερομηνία τροποποίησης είναι προηγούμενη από τις εξαρτήσεις του (δηλαδή foo. ο και bar. ο) αυτό θα ξαναχτιστεί. Εάν κάποια από εξαρτήσεις απουσιάζει η ίδια λογική ισχύει και για αυτή. Με αυτό τον τρόπο εάν χτίσουμε μια φορά το foobar, και μετά τροποποιήσουμε το bar. c και ξανατρέξουμε make foobar, το make θα θεωρήσει αναδρομικά ότι:

- το bar. ο είναι out of date όσον αφορά την εξάρτηση bar. c
- Όταν bar. ο έχει πλέον μια πιο πρόσφατη ημερομηνία μετατροπής από το foobar και για αυτό το τελευταίο είναι out of date όσον αφορά την dependencyτου bar. ο, έτσι χρειάζεται να ξανχτιστεί.

με αυτόν τον τρόπο το make πετυχαίνει μια σχεδόν βέλτιστη στρατηγική για την επίτευξη κάθε φορά του ελάχιστου ποσοστού των απαιτούμενων στόχων. Τώρα που ξεκαθαρίσαμε την βασική λογική των make ας κάνουμε πιο σαφή μερικά από τα βασικά χαρακτηριστικά τους που κάνουν τη ζωή μας πιο εύκολη.

## (a) Phony targets

Μερικές εργασίες δεν είναι αρχεία και χρειάζονται να τρέχουν κάθε φορά που το make τις συμπεριλαμβάνει στο dependency tree. Γι αυτά έχουμε ένα ειδικό keyword .PHONY:.

Παρακάτω είναι ένα παράδειγμα.

```
.PHONY:
clean:
    rm -rf *
```

Αυτό λέει στο make ότι κανένα αρχείο ονομαζόμενο clean δε θα δημιουργηθεί τρέχοντας rm -rf \*, και επίσης ακόμα και εάν υπάρχει ένα up-to-date ονομαζόμενο αρχείο ονομαζόμενο clean, αυτό το target θα τρέχει ανεξάρτητα.

Αξίζει να σημειώσουμε ότι οι phony εξαρτήσεις πάντα θα θεωρούνται out of date.

Για παράδειγμα:

```
.PHONY:
say-hello:
    echo "hello"

test.txt: say-hello
    touch test.txt
```

Όταν το touch test.txt θα τρέχει κάθε φορά που τρέχουμε make test.txt απλώς γιατί το make δεν μπορεί να γνωρίζει με βεβαιότητα ότι το phony target say-hello δεν άλλαξε τίποτε σημαντικό για το test.txt. Για αυτό το λόγο τα phony targets χρησιμοποιούνται για user facing tasks.

#### (b) Variables

Τα makefiles μπορούν να έχουν μεταβλητές ορισμένες με ποιχίλους τρόπους. Μερικές περιπτώσεις που έχουν γίνει για να χρησιμοποιηθούν στην wikiepedia-mirror παρουσιάζονται παρακάτω.

Αναδρομικές μεταβλητές

```
OBJCETS = foo.o bar.o
show:
    echo $(OBJECTS)
```

Τρέχοντας make show θα εμφανίσει foo.o bar.o στην πονσόλα. Όλες οι μεταβλητές αντιπαθιστώνται με τις τιμές τους αν βάλει πάνεις πρενθέσεις γύρω από το όνομά τους παι προθέσει ένα δολάριο (\$). Οι μεταβλητέσ των makefiles δεν έχουν τύπους, αναφορά σε μια μεταβλητή είναι ισοδύναμη με string substitution, όπως είναι παι στο shell scripting.

Οι μεταβλητές που ορίζονται με ένα απλό \= είναι recursively expanded. Αυτό σημαίνει ότι αφού το όνομα της μεταβλητής αντικαθίσταται από την τιμή της μια αναδρομική διαδικασία συνεχίζει να κάνει expand τις τιμές που προκύπτουν με την ίδια τη μεταβλητή ακόμα στο local scope.

```
library = foo

foo-libs = -lfoo
foo-includes = -I./include/foo

bar-libs = -lbar
bar-includes = -I./include/bar

libs = $($(library)-libs)
includes = $($(library)-includes)

waz:
    gcc waz.c $(includes) $(libs)

τοέχοντας make

gcc waz.c $(includes) $(libs)

gcc waz.c $($(library)-includes) $($(library)-libs)

gcc waz.c $($(library)-includes) $($(library)-libs)

gcc waz.c $(foo-includes) $(foo-libs)

gcc waz.c -I./include/foo -lfoo
```

Παρατηρήστε πως οι αναφορές στις μεταβλητές καθαυτές δημιουργήθηκαν.

Μεταβλητές μπορούν επίσης να ορισθούν στην εντολή make

```
$ make --just-print library=bar
gcc waz.c -I./include/bar -lbar
```

## ii. Simple variables

Μερικές φορές δεν είναι επιθυμητό για τις μεταβλητές να είναι expanded επ αόριστον:

```
kurma = the world $(support1)
animal1 = four elephants
animal2 = tortoise
support1 = supported by $(animal1) $(
    support2)
support2 := supported by a $(animal2) $(
    support2)
all:
    echo $(kurma)
```

Εδώ πος σπαθουμε να δημιουργήσουμε ένα άπειρο μήνυμα.

```
$ make --just-print
Makefile:5: *** Recursive variable `support2
   ' references itself (eventually). Stop.
```

το σύστημα μεταβλητών δηλαδή είναι κατά κάποιον τοόπο total[26], με άλλα λογία η εύρεση της τιμής μεταβλητών μπορεί να είναι αναδρομική άλλα πρέπει να τερματίζει. Μπορούμε να αποφύγουμε αυτόν τον περιορισμό ορίζοντας μεταβλητές με :=:

```
make --just-print
echo the world supported by four elephants
    supported by a tortoise
```

#### iii. Automatic variables

Το Makefile επίσης ορίζει μερικές contextual μεταβλητές οι οποίες είναι ορισμένες. Οι πιο σημαντικές automatic variables που ορίζει το gnu make είναι οι ακόλουθες

- \$0: Το όνομα του αρχείου του target. Εάν το target είναι ένα archive member, τότε \$0 είναι το όνομα του archive αρχείου. Στο pattern rule που έχει πολλαπλά targets, \$@ είναι το όνομα του οποιουδήποτε target που κάνει το rule's recipe να τρέχει.
- \$%: Το όνομα τουtarget member, όταν το target είναι ένα archive member. Για παράδειγμα, εάν το target είναι foo.a(bar.o) τότε \$% είναι bar.o και \$@ είναι foo.a. \$% είναι άδειο όταν το target δεν είναι ένα archive member.
- \$<: Το όνομα του πρώτου prerequisite. Εάν το target πήρε το recipe του από έναν implicit rule, αυτό θα είναι το πρώτο prerequisite που προστέθηκε από το implicit rule.
- \$?: Τα ονόματα από όλες τις εξαφτήσεις που είναι νεότεφα από το target, με μενά μεταξύ τους . Για τα prerequisites που είναι archive members, μόνο named member χφησιμοποιούνται (βλέπε Archives).
- \$^: Τα ονόματα όλων των prerequisites, με κενά μεταξύ τους. Για τα prerequisites τα οποία είναι archive members,

μόνο των named member χοησιμοποιείται. ένα target έχει μόνο ένα prerequisite σε κάθε άλλο αρχείο από το οποίο εξαρτάται, αναξαρτήτως από το πόσες φορές κάθε αρχείο είναι καταχωρημένο ως ένα no matter how many times each file prerequisite. Έτσι εάν τοποθετήσουμε στη λίστα ένα prerequisite για περισσότερο από μια φορά για ένα target, η value του \$^ περιέχει μόνο ένα αντίγραφο του ονόματος.

## iv. Συναρτήσεις

Οι συναρτήσεις είναι παρόμοιες με μεταβλητές ως προς το ότι και αυτές γίνονται expand σε συμβολοσειρές. Η μόνη διαφορά είναι ότι επιδέχονται παραμέτρους.

## 3. Πηγαίοι κωδικές

(a) page<sub>remover.c</sub>

```
/*
 * Copyright 2014 Chris Perivolaropoulos <
        cperivol@csail.mit.edu>
 *
 * This program is free software: you can
        redistribute it and/or
 * modify it under the terms of the GNU General
        Public License as
 * published by the Free Software Foundation,
        either version 3 of the
 * License, or (at your option) any later
        version.
 *
```

```
st This program is distributed in the hope that
    it will be useful, but
 * WITHOUT ANY WARRANTY; without even the
    implied warranty of
 * MERCHANTABILITY or FITNESS FOR A PARTICULAR
    PURPOSE.
 * See the GNU General Public License for more
    details. You should
 * have received a copy of the GNU General
    Public License along with
 * this program.
 * If not, see <a href="http://www.gnu.org/licenses/">http://www.gnu.org/licenses/>.
 * This should fill a range in a file with
    spaces. This is an in-place
 * operation so it should be pretty fast.
 * Usage: page_remover PATH OFFSET LENGHT
#include <assert.h>
#include <fcntl.h>
#include <pthread.h>
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <sys/mman.h>
#include <sys/stat.h>
#include <sys/types.h>
#include <semaphore.h>
#include <unistd.h>
#include <unistd.h>
#define USAGE_INFO "page_remover_PATH_OFFSET_
#define PRINT(ctx, args...) do { sem_wait(&ctx->
   stdio_mutex); \
                                 printf(args);
                                 fflush(stdout);
                                 sem_post(&ctx->
```

```
stdio_mutex);
                                } while(0)
typedef struct context {
    int fd;
    size_t size;
    off_t off;
    sem_t stdio_mutex;
    void* data;
} context_t;
context_t* context_init(char* fname, off_t off,
   size_t len)
{
    context_t * ctx = (context_t*)malloc(sizeof(
       context_t));
    off_t pa_off = off & ~(sysconf(_SC_PAGE_SIZE
       ) - 1);
    sem_init(&ctx->stdio_mutex, 0 /* Shared.
       Usually ignored */ , 1);
    PRINT(ctx, "Opening_\%s_at_\%lu_(len:_\%lu)\n",
        fname, off, len);
    ctx->off = off-pa_off;
    ctx->fd = open(fname, O_RDWR, 0x0666);
    if (ctx->fd == -1) {
      perror("open");
      return NULL;
    }
    ctx->size = len;
    ctx->data = mmap(0, len+ctx->off, PROT_READ
       | PROT_WRITE,
                   MAP_SHARED, ctx->fd, pa_off);
    if (ctx->data == MAP_FAILED) {
      perror ("mmap");
      return NULL;
    return ctx;
}
```

```
void context_destroy(context_t* ctx)
         if (close (ctx->fd) == -1)
           perror ("close");
         if (munmap ((void*)ctx->data, ctx->size) ==
           perror ("munmap");
         sem_destroy(&ctx->stdio_mutex);
         free(ctx);
     }
     int main(int argc, char *argv[])
         if (argc != 4)
           fprintf(stderr, USAGE_INFO);
         context_t *ctx = context_init(argv[1], atoi(
            argv[2]), atoi(argv[3]));
         /* \ \textit{You MIGHT want to thread this but I dont}
             think it will make
          * much more difference than memset. */
         memset(ctx->data + ctx->off, 'u', ctx->size)
         context_destroy(ctx);
         return 0;
    }
(b) utf8thread.c
     #include <assert.h>
     #include <fcntl.h>
     #include <pthread.h>
     #include <stdio.h>
     #include <stdlib.h>
    #include <string.h>
    #include <sys/mman.h>
     #include <sys/stat.h>
     #include <sys/types.h>
     #include <semaphore.h>
     #include <unistd.h>
```

```
#include <unistd.h>
sem_t stdio_mutex;
#define PRINT(args...) do {sem_wait(&stdio_mutex
      printf(args);
      fflush(stdout);
      sem_post(&stdio_mutex);
    } while(0)
/* #define DEBUG(args...)
                                        PRINT (
   args) */
#define DEBUG(...)
#define DEFAULT_CHAR '_'
#define WORKERS 8
#define MESSAGE_DENSITY 100000000
typedef unsigned long long u64;
#define UTF_LC(1) ((0xff >> (8 - (1))) << (8 - (
   1)))
#define UTF_CHECK(1, c) (((UTF_LC(1) & (c)) ==
   UTF_LC(1)) && (0 == ((c) & (1 << (7-(1)))))
#define UTF_LEN(x) (UTF_CHECK(6, x) ? 6 :
   \
                  UTF_CHECK(5, x) ? 5 : \
                  UTF_CHECK(4, x) ? 4 : \
                  UTF_CHECK(3, x) ? 3 : \
                  UTF_CHECK(2, x) ? 2 : -1)
struct crange {
    u64 start, end;
};
/* Get return the next character after the last
   correct one. */
inline u64 valid_utf8(u64 c)
```

```
char i;
    /* Ascii */
    if ((*(char*)c & 0x80) == 0)
      return c+1;
    /* */
    for (i = UTF_LEN(*(char*)c)-1; i>0; i--) {
      if (!UTF_CHECK(1, *(char*)c)) {
          return (u64) NULL;
      }
    }
    return i<0 ? 0 : c+1;
}
void* fix_range(void* _r)
    struct crange* r = _r;
    u64 tmp, id = r \rightarrow start;
    long long unsigned count = 0;
    while ((u64)r\rightarrow start < (u64)r\rightarrow end) {
      if (count++ % MESSAGE_DENSITY == 0)
          printf ("[worker: \( \)0x\%016llx] \( \)Done \( \)with
              if (!(tmp = valid_utf8(r->start))){
          PRINT("Invalid_char_0x%x_(next:_0x%x)\
              n",
                *(char*)r->start, *(char*)(r->
                    start+1));
          *((char*)r->start) = DEFAULT_CHAR;
          (r->start)++;
      } else {
          r->start = tmp;
      }
    }
    PRINT ("[worker: \_0x\%016llx]\_0UT\n", id);
    return NULL;
}
void run(u64 p, u64 sz)
```

```
{
    int n, i;
    u64 wsize;
    pthread_t workers[WORKERS];
    struct crange rngs[WORKERS];
    wsize = sz/WORKERS + 1;
    printf("Base_address:_0x%016llx,_step_size:_
       0x\%016llx\n", p, wsize);
    for (i=0; i<WORKERS; i++){</pre>
      rngs[i].start = p + wsize*i;
      rngs[i].end = p + wsize*i + wsize;
      PRINT("Spawning_worker_%d_on_range_[0x%016]
          llx, 0x\%016llx), %llu bytes..., i,
          rngs[i].start, rngs[i].end, wsize);
      if ((n = pthread_create(workers+i, NULL,
          fix_range, (void*)(rngs+i)))) {
          PRINT("FAIL\n");
          perror("worker");
          return;
      }
      PRINT("OK\n");
    }
    PRINT ("Wrapping up...\n");
    for (i=0; i<WORKERS; i++) {</pre>
      PRINT ("Joining worker \%d...", i);
      pthread_join(workers[i], NULL);
      PRINT ("OK\n");
      PRINT("Worker_\%d_went_\through_\%llu_bytes.\
            i, (u64)rngs[i].end - (u64)rngs[i].
                start);
    }
}
int main(int argc, char *argv[])
{
    int fd;
    long long int sz, p;
    struct stat buf;
```

```
ignored */, 1);
         fd = open(argv[1], O_RDWR, 0x0666);
         if (fd == -1) {
           perror("open");
           return 1;
         fstat(fd, &buf);
         sz = buf.st_size;
        printf("File_size:_0x%016llx\n", sz);
         p = (u64)mmap (0, buf.st_size, PROT_READ |
            PROT_WRITE , MAP_SHARED, fd, 0);
         if (p == -1) {
          perror ("mmap");
           return 1;
        run(p, buf.st_size);
         if (close (fd) == -1) {
          perror ("close");
           return 1;
         if (munmap ((void*)p, buf.st_size) == -1) {
           perror ("munmap");
           return 1;
        }
         sem_destroy(&stdio_mutex);
        return 0;
    }
(c) sql-clear.sh
    #!/bin/bash
    MUSER="$1"
    MPASS="$2"
    MDB="$3"
```

sem\_init(&stdio\_mutex, 0 /\* Shared. Usually

```
# Detect paths
       AWK=$(which awk)
       GREP=$(which grep)
       if [ $# -ne 4 ]
       then
                echo "Usage: \( \sqrt{SQL-User-Name} \) \( \lambda \) \( \text{MySQL-} \)
                     User-Password}_{\| \{ MySQL-Database-Name}_{\| \{ \}
                     MySQL_uexecutable_uto_use"
                echo "Drops_{\sqcup}all_{\sqcup}tables_{\sqcup}from_{\sqcup}a_{\sqcup}MySQL"
                exit 1
       fi
       TABLES=$($MYSQL -u $MUSER -p$MPASS $MDB -e 'show
             tables' | $AWK '{ print $1}' | $GREP -v '^
            Tables')
       for t in $TABLES
       do
                \textbf{echo} \ \texttt{"Clearing} \bot \$t \bot table \bot from \bot \$MDB \bot database
                $MYSQL -u $MUSER -p$MPASS $MDB -e "
                     truncate_table_$t"
       done
(d) webmonitor.py
       .....
       Just_{\sqcup}feed_{\sqcup}pairs_{\sqcup}of
       <epocudate>u<floatuvalue>
       or_{\sqcup}even_{\sqcup}just
       <float uvalue>
       {\tt One} {\tt \sqcup} {\tt way} {\tt \sqcup} {\tt to} {\tt \sqcup} {\tt do} {\tt \sqcup} {\tt that} {\tt \sqcup} {\tt would} {\tt \sqcup} {\tt be}
       \squarepython\squarewebmonitor.py
       \verb"and" \sqcup \verb"will" \verb"plot" \verb"them" \verb"on" \verb"port" \verb"8888. "This" \verb"will" \verb"
```

MYSQL=\$4

```
\verb|also|| \verb|pipe|| the|| \verb|input|| right|
out_{\sqcup}to_{\sqcup}the_{\sqcup}output._{\sqcup}Strange_{\sqcup}input_{\sqcup}will_{\sqcup}be_{\sqcup}ignored
    \sqcupand\sqcuppiped\sqcupthis\sqcupway,
but_{\sqcup}this_{\sqcup}needs_{\sqcup}to_{\sqcup}be_{\sqcup}done_{\sqcup}by_{\sqcup}awk_{\sqcup}aswell_{\sqcup}in_{\sqcup}the_{\sqcup}
    above \_ example.
.....
import sys
import json
import time
from threading import Thread
from collections import deque
import tornado.websocket as websocket
import tornado.ioloop
import tornado.web
HTML = """
<!DOCTYPE_HTML_PUBLIC_"-//W3C//DTD_HTML_4.01//EN</pre>
     "u"http://www.w3.org/TR/html4/strict.dtd">
<html>
⊔⊔<head>
\square\square\square\square <meta\squarehttp-equiv="Content-Type"\squarecontent="
    text/html; charset=utf-8">
⊔⊔⊔⊔<title>DrNinjaBatmans Uebsockets</title>
uuuu <script utype="text/javascript" usrc="http://
     code.jquery.com/jquery-1.10.1.js"></script>
uuuu <script utype="text/javascript" usrc="http://
    code.highcharts.com/highcharts.js"></script>
\Box\Box\Box\Box < script >
var uchart; u//uglobal
var_{\sqcup}url_{\sqcup} = location.hostname_{\sqcup} + l_{\sqcup}': '_{\sqcup} + l_{\sqcup}(parseInt(
    location.port));
var_ws_=unew_WebSocket('ws://'u+urlu+u'/
    websocket');
ws.onmessage\square = \squarefunction(msg)\square{
⊔⊔⊔⊔add_point(msg.data);
};
//_{\sqcup}ws.onclose_{\sqcup}=_{\sqcup}function()_{\sqcup}{_{\sqcup}alert('Connection_{\sqcup}
    closed.');<sub>||</sub>};
```

```
var_{\sqcup}add_{point_{\sqcup}}=_{\sqcup}function(point)_{\sqcup}
uuuuuvaruseriesu=uchart.series[0],
\sqcup \sqcup \sqcup \sqcup \sqcup \sqcup \sqcup shift\sqcup = \sqcup series.data.length<math>\sqcup > \sqcup \%d;
\Box\Box\Box\Box\Box chart.series[0].addPoint(eval(point),\Boxtrue,\Box
    shift);
};
$(document).ready(function(),{
\sqcup \sqcup \sqcup \sqcup \sqcup chart \sqcup = \sqcup new \sqcup Highcharts. Chart (JSON.parse ('%s
     '));
});
\Box\Box\Box\Box </script>
uu</head><body><divuid="container"ustyle="width:
    \square 800px; \square height: \square 400px; \square margin: \square 0 \square auto" > </div
    ></body></html>
config = {
      'visible_points': 10,
      'py_chart_opts': { 'chart': { 'renderTo': '
          container',
                                                     defaultSeriesType
                                                     ': 'spline
                                                     '},
                                 'title': { 'text': '
                                     DrNinjaBatmans ⊔ data'
                                     },
                                 'xAxis': { 'type': '
                                     datetime',
                                                     tickPixelInterval
                                                     ': '150'},
                                 'yAxis': { 'minPadding':
                                     0.2,
                                                'maxPadding':
                                                     0.2,
                                                'title': {'
                                                    text': '
                                                     Value',
                                                                   margin
```

```
80}
                               },
                        'series': [{ 'name': '
                           Data',
                                      'data':
                                         []}]}
}
def date_float(s):
    try:
        date, val = s.split()
    except ValueError:
        val = s.strip()
        date = time.time()
    return int(date), float(val)
def send_stdin(fn=date_float):
    for raw in sys.stdin:
        sys.stdout.write(raw)
        # Ignore strange input.
        try:
            jsn = json.dumps(fn(raw))
            buf.append(jsn)
            for w in websockets:
                try:
                    w.write_message(jsn)
                except websocket.
                    WebSocketClosedError:
                    pass
        except:
            pass
    for ws in websockets:
        ws.close()
```

class StdinSocket(websocket.WebSocketHandler):

```
def open(self):
             for i in buf:
                 self.write_message(i)
             websockets.append(self)
        def closs(self):
             websockets.remove(self)
    class MainHandler(tornado.web.RequestHandler):
         def get(self):
             self.write(HTML % (int(config['
                visible_points']),
                                json.dumps(config['
                                    py_chart_opts']))
    if __name__ == "__main__":
         application = tornado.web.Application([
             (r"/", MainHandler),
             (r'/websocket', StdinSocket),
        ])
        buf = deque(maxlen=int(config['
            visible_points']))
        websockets = []
        config['args'] = []
        for a in sys.argv[1:]:
             if '=' in a:
                 k, v = a.split('=', 1)
                 config[k] = v
                 config['args'].append(a)
        Thread(target=send_stdin).start()
         application.listen(8888)
        tornado.ioloop.IOLoop.instance().start()
(e) xml-parse.sh
```

#!/bin/bash

```
# Simply removing specific articles fixes the
   xerces error with
# UTF8. If the articles are alone the error goes
     away
# aswell. Extremely weird but that's life.
   Fortunately the article is
# just a stub about some toad (Cranopsis
   bocourti)
# xml-parse.sh ORIGINAL_XML
    TITLE_OF_ARTICLE_TO_REMOVE [inplace]
# if `inplace` is there the c program will be
    used to cover the article
# with spaces. This is much faster. Should be
    anyway. Otherwise the
\# page is just ommitted and the result is dumped
    in stdout. Helping
# messages are dumped in stderr After this you
   can run:
# java -jar tools/mwdumper.jar RESULTING_XML --
   format = sql: 1.5 > SQL_DUMP
set -e
set -o pipefail
if [[ $# -lt 2 ]]; then
    echo "xml-parse.sh_{\sqcup}ORIGINAL_{\bot}XML_{\sqcup}
        TITLE_OF_ARTICLE_TO_REMOVE_[inplace]"
        1>&2
    exit 0
fi
function my_dd {
    coreutils_version=$(dd --version | head -1 |
         cut -d\ -f3 | colrm 2 2 )
    if [[ $coreutils_version -ge 822 ]]; then
      eval "dd_{\sqcup}iflag = count_bytes_{\sqcup}iflag = direct_{\sqcup}
          oflag=seek_bytes_{\perp}ibs=1M_{\perp}$0"
    else
      echo "Your_coreutils_may_be_a_bit_old_(
          $coreutils_version). □822 □ is □ the □ one □
          cool_{\perp}kids_{\perp}use." > &2
```

```
eval "dd_{\sqcup}$0_{\sqcup}ibs=1"
     fi
}
ORIGINAL_XML=$1
# Dump a part of the file in sdout using dd.
# Usage:
# file_range <filename> <first_byte> <start/end/</pre>
    length>
# Length can be negative
function file_range {
     file=$1
     start=$2
     len=$3
     case $len in
       "end") my_dd if=$file skip=$start || exit
           1; return 0;;
       "start") my_dd if=$file count=$start ||
           exit 1; return 0;;
       "") echo "len_{\sqcup}was_{\sqcup}empty_{\sqcup}(file:_{\sqcup}$file,_{\sqcup}
           start: \( \start \, \( \len \square \start \). \( \len \correct \)
           format_{\sqcup} < filename >_{\sqcup} < byte_{\sqcup} start >_{\sqcup} < length
           |'start'|'end'>" 1>&2; exit 1;;
       *);;
     esac
     if [[ $len -gt 0 ]]; then
       # Dump to stdout
       my_dd if=$file skip=$start count=$len ||
           exit 1
     else
       skip=$(($start + ($len)))
       len=$((- ($len)))
       if [[ skip - 1t 0 ]]; then
            skip=0
            len=$start
       fi
       # Dump to stdout
         my_dd if=$file skip=$skip count=$len ||
```

```
exit 1
    fi
}
function backwards {
    tac -b | rev
function byte_offset {
    grep -b -o -m 1 -F "$1" | cut -d : -f1
}
# Throw everything but the page in stdout
# neg_xml_page "Barack Obama"
function neg_xml_page {
    term="<title>$1</title>"
    title_offset=$(cat $ORIGINAL_XML |
        byte_offset "$term")
    echo -e "\n\tMethod: \subseteq$2(blank \subseteq is \subsete ok) " 1>&2
    echo -e "\tsearch_term:_\$term" 1>&2
    echo -e "\tfile:⊔$ORIGINAL_XML" 1>&2
    echo -e "\ttitle offset: stitle offset" 1>&2
    # Fail the term is invalid
    if [ -z "$title_offset" ]; then
       echo "Found" '$title_offset' Grep-ing (cat )
          \verb| \_\$ORIGINAL_XML| | | grep| - b| - m| 1| - F| \setminus "
          term''_{\square}|_{\square}cut_{\square}-d:_{\square}-f1)" 1>&2
       exit 1
    fi
    to_page_start=$(($(file_range $ORIGINAL_XML
        $title_offset -1000 | backwards |
        byte_offset "(echo_{\sqcup}'<page>'_{\sqcup}|_{\sqcup}rev)")+7)
    echo -e "\tto⊔page⊔start⊔(relative):⊔
        $to_page_start" 1>&2
    file_range $ORIGINAL_XML $title_offset end |
         byte_offset "</page>" >&2
    echo $(($(file_range $ORIGINAL_XML
        $title_offset end | byte_offset "</page>
        ")+7)) >&2
    to_page_end=$(($(file_range $ORIGINAL_XML
```

```
$title_offset end | byte_offset "</page>
        ")+7)) # len('</page>') == 7
    echo -e "\tto⊔page⊔end⊔(relative):⊔
        $to_page_end" 1>&2
    page_start=$(($title_offset - $to_page_start
         +1 ))
    echo -e "\tpage_start: \_$page_start" 1>&2
    page_end=$(($title_offset + $to_page_end))
    echo -e "tpage_lend:_l$page_end" 1>&2
    echo -e "\tbytes_to_copy:_\$(($(du_-b_
        SORIGINAL_XML_{\square}|_{\square}cut_{\square}-f1)_{\square}-_{\square}page_start_{\square}+
        □$page_end))" 1>&2
    echo "Going Lto Lcopy L$page_start bytes" 1>&2
    file_range $ORIGINAL_XML $page_start start
    echo "Finished the first half up to
        page_start, u$((u$(du_b)$ORIGINAL_XML_l)
        \lfloor \text{cut} \rfloor - f \rfloor 1) \rfloor - \lfloor \text{spage\_end} \rfloor) \rfloor to \rfloor go = 1 \times 2
    file_range $ORIGINAL_XML $page_end end
    echo "Finished \perp the \perp whole \perp thing." 1>&2
}
# Put stdin betwinn mediawiki tags and into
    stdout
function mediawiki_xml {
     (head -1 $ORIGINAL_XML; sed -n "/<siteinfo
        >/,/<\/siteinfo>/p;/<\/siteinfo>/q"
        $ORIGINAL_XML ; cat - ; tail -1
        $ORIGINAL_XML )
}
# 1: XML File
# 2: Article
# 3: Method (leave blank)
# Assert that the file is there and is not empty
fsize=$(du -b $ORIGINAL_XML | cut -f1)
if [[ 0 -eq $fsize ]]; then
     echo "ERROR: _empty_xml_file_$ORIGINAL_XML"
        1>&2
    exit 1
fi
```

```
echo "Will_remove_article_'\$2'_from_file_\$1_(
    size:_\$fsize)" 1>&2
if ! neg_xml_page "$2" "$3"; then
    ret=$?
    echo "XML_parsing_script_failed" 1>&2
    exit $ret;
fi
```

# Part IV Βιβλιογοαφία