# Robust Morphologic Analyzer
# for highly inflected Languages

## Andrés Tomás Hohendahl [1,3]

## José Francisco Zelasco [1,2]

1 Lab. de Estereología y Mecánica Inteligente
  Facultad de Ingeniería, U.B.A.

2 INTIA Facultad de Ciencias Exactas, UNCPBA

3 Instituto de Ingeniería BioMédica, U.B.A.

# Languages & Dictionaries

- **Inflected Languages**
  - Slighltly Inflected:
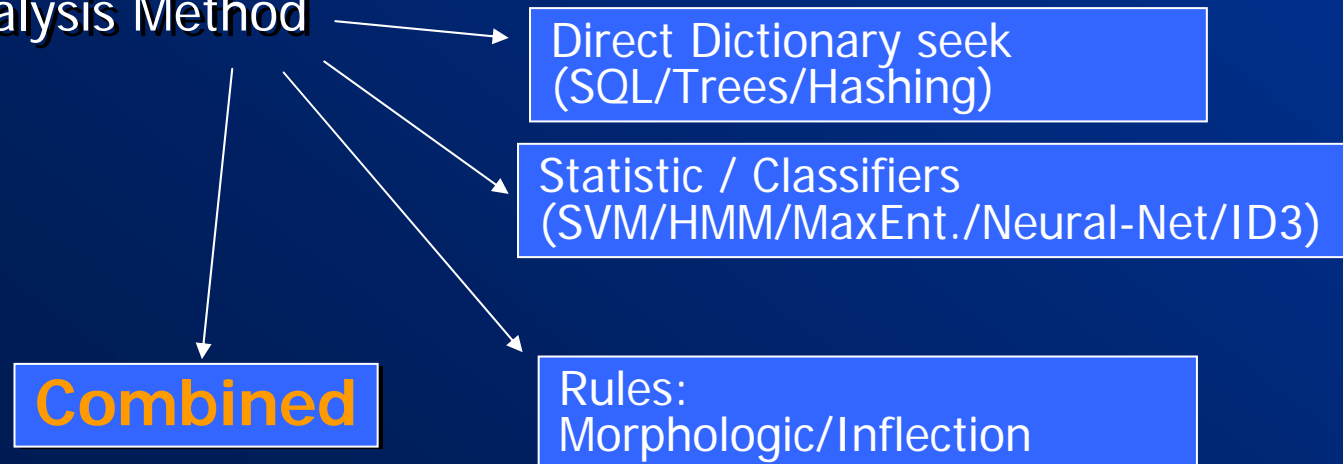    - English ~ 80k roots (x2.3) < 200k total
  - Highly Inflected + Parasynthetic
    - Spanish ~ 80k roots (x10k) > 3000M total
    - Huge Word-Space 1.0 E22 words for 15 letters
    - Similar for Polish, French, Italian, Portuguese, etc.
  - Analysis Method

Direct Dictionary seek
(SQL/Trees/Hashing)

Statistic / Classifiers
(SVM/HMM/MaxEnt./Neural-Net/ID3)

**Combined**

Rules:
Morphologic/Inflection

# Languages & Dictionaries

- Lexical Word analysis (goals)
  - Minimum stored amount of data
  - Obtain Semantic and Grammatical Information
  - Tolerate Misspelling & Suggest Corrections
  - Do all above: efficiently

- Used Method
  - Store word-Roots along with applicable Rules & base Tags
  - Each Morphologic Rule contains Grammatical & Semantic info.
  - Fast in-memory data structures: (Patricia-Trie y Tst)
  - Recursive & Greedy Algorithm: Seek / Lemmatize
  - Intelligent Spelling Suggestion: min. seek / max. probability
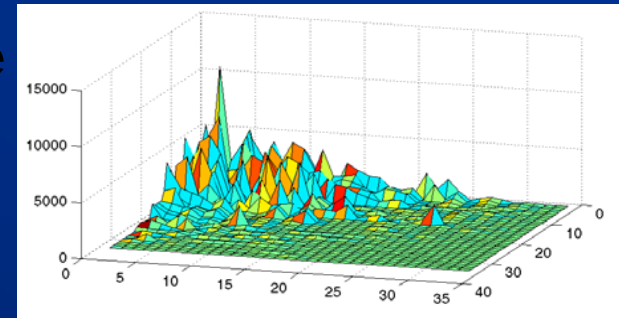
# Languages & Dictionaries

- ## Language Recognition

  ### Statistic

  - Fast (few operations)
  - Compact Datasets ~5kb/language
  - Good Recall (F-Score >95%)
  - Reduces Unnecessary Seeks
  - State of the Art in 2004

  ### Other Methods
  Proprietary (MS, etc.)
  Brute Force (high cost)

- Hohendahl, A.T. Zelasco, J.F. WICC 2006 (art.694)



ES- Distribución de diletras
por: inicial, segunda (48424words)

- Padró, Lluís/Munsa. TALP 2004 UPC

# Efficient Index Structures

- SQL / TSQL / CQL (high level)
  - → Inefficient for partial words
  - → High cost (resources, licensing, maintenance, TOC)
- Binary Trees, M-Trees, Radix-Trees
  - → Less efficient for partial-matching

✓ Tries & TST (Ternary Search Trees)
  - → Linear Time **O**(word length)
  - → Useful for error detection/correction
  - → Easy to finding Sub-Ranges for Similarity
  - → Flexible: Linkable & Combinable

# Reversible Morphologic Rules

- **Spanish** vocabulary using enhanced ASPELL compression
  - ~ 3.900 inflection rules (300 prefix/infix + 3600 suffix)
  - ~ 200 Semantic/Grammatical Attributes.
  - ~ 79.000 Root words (Lemmas)
  - ~ 300 kb compacted (*.zip)

ASPELL.org (GNU)

Yields ➜ >5 Million exact recognizable words
+ Phonetic Guess Sampa (Sound-Like)
+ Enhanced Spell Correction (statistic-guess)
+ Morphologic Guessing (statistic + rules)
+ Parasynthesis (multiple combinations)

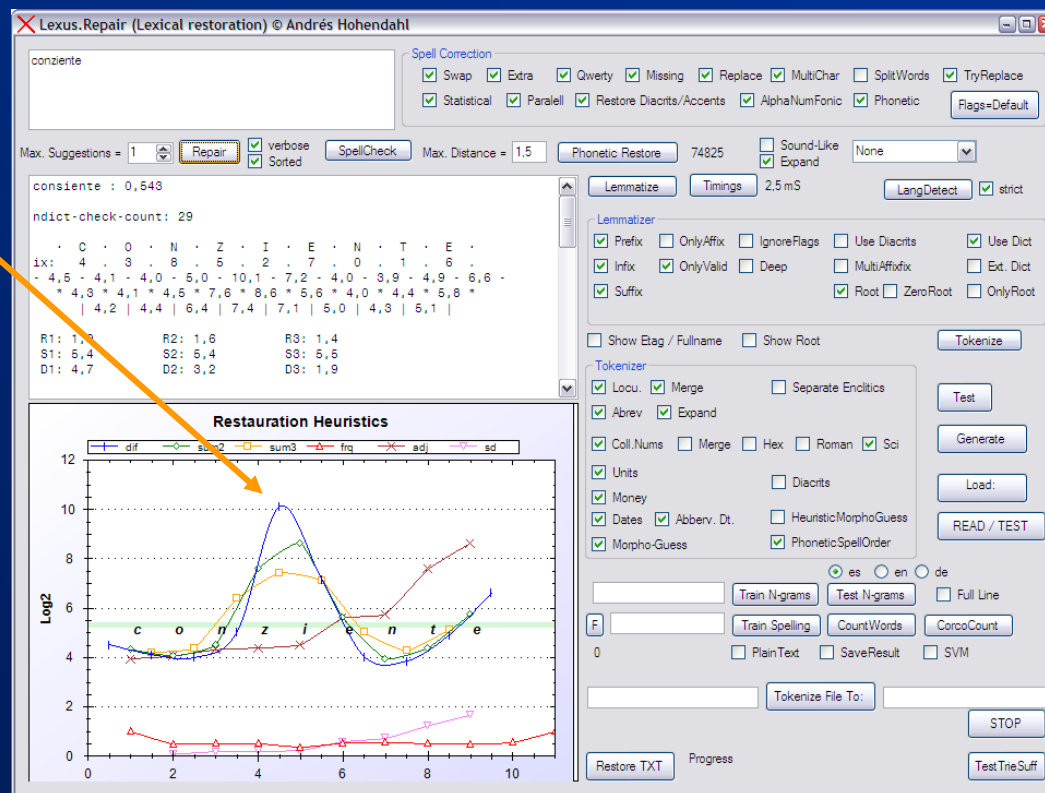➜ Huge >5000 Million word space !
(not including correctable mis-spells)

# Morphologic Analysis Algorithm

- Pseudo-Logic Diagram (very simplified) for Finding a Word

*if (word in Roots) → found*

*Acumulate word in [seePrefix]*

*foreach Affix Rule in Suffix-Rules*

*if Rule Applicable to word → strip-Suffix*

*if (stripped in Roots & bears Rule)*
*→ found*

*else accumulate in [seePrefix]*


*foreach word in [seePrefix]*

*foreach Affix Rule in Prefix-Rules*

*if Rule Applicable to word → strip-Prefix*

*if (stripped in Roots & bears Rule)*
*→ found*

# Spell Error Detection-Correction

- Using Bigram & Trigram Freq. from Language Detector.
- Heuristics to find best fitting replacement.
- Reduced seek count.
- Detects promptly unusual zones.
- Uses language specific rules.
- Usually finds the best "human" word in the first (few) trials
- Shares TST/Tries with Analyzer.
- Simply based on: **Poor-Man-Speller**

# Dictionary + Rule Editor/Utility

## Runs on Windows (.NET 2.0 C# Platform)

**Features:**

**Builds Rules**

**Test Rules**

**Analyzes words**

**Expands words**

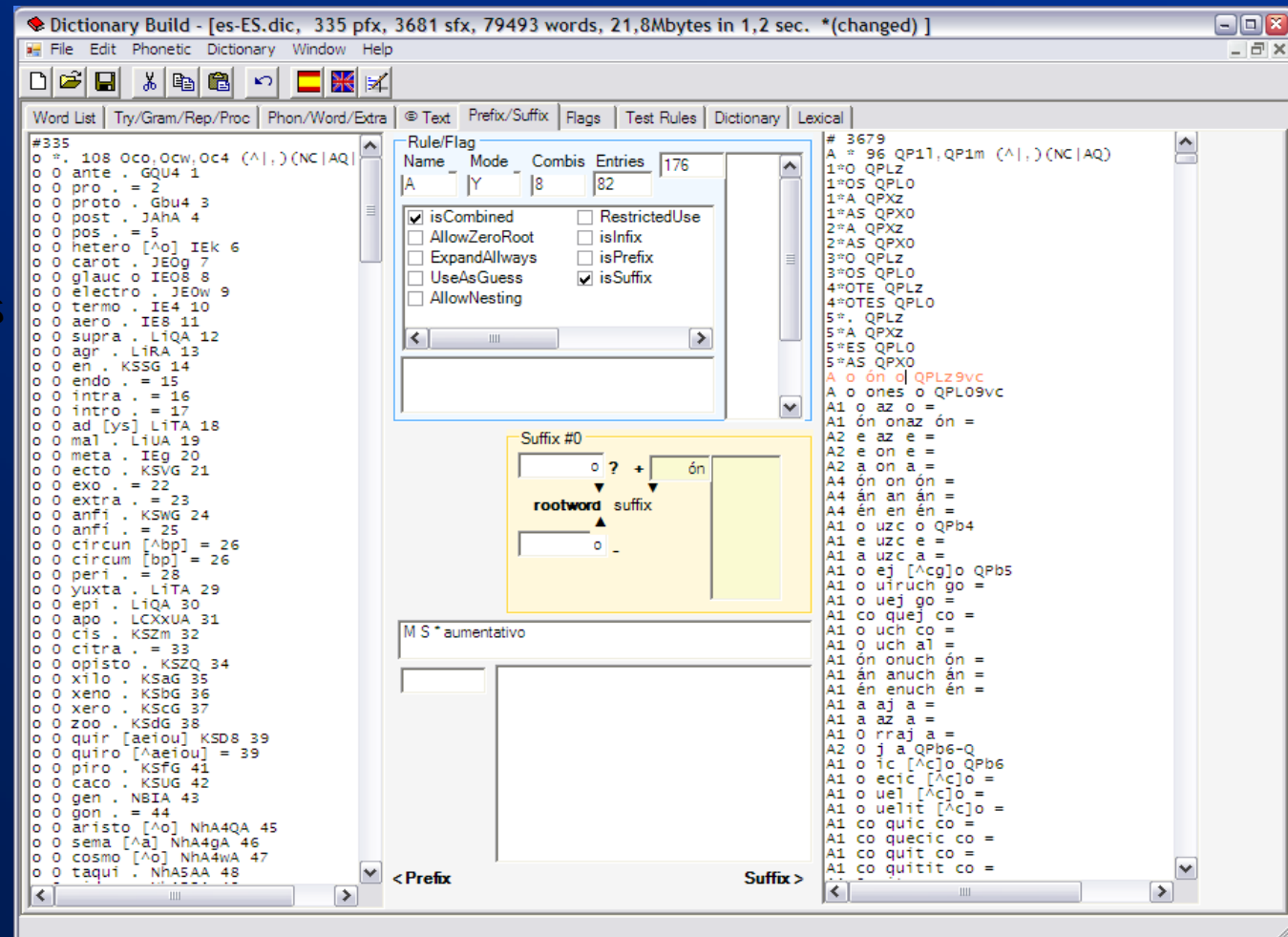**Benchmarks**

**Imports**

   IFFIX

   ASPELL

   Word-List

**Does all kind of word/tag operations**

# Phonetic Similarity Module

**New Algorithm**

– Measures what humans "think-sounds" written text.

– Based on analogical phono-articulartory model
(uses non-linear kernel on: pitch-vibration, nasal-lingual position/occlusion + openness + fricative energy )

– Measures Distance among words with a Real Number 0..max where (1.0 is the mean phonemic distance)

– Correlation with human-perception over 85%.

– Establishes a good parameter for spell correction delivering the correct word even with worse misspells,

Example:

***VAHIEMA → BALLENA***

(only 1 guess d~0.69 0.001sec)

A. Hohendahl,
S. Zanutto,
A.Wainselboim
SLAN 2007

# Phonetic Similarity Algorithm

## Highlights

– Very little literature found on the subject.

– Outperforms classic Lexical distance for cognitive perception experiments and measures (Levenshtein, etc.)

– Very Fast (over 30k/second)

– Small Memory usage

**Windows Utility for Testing**
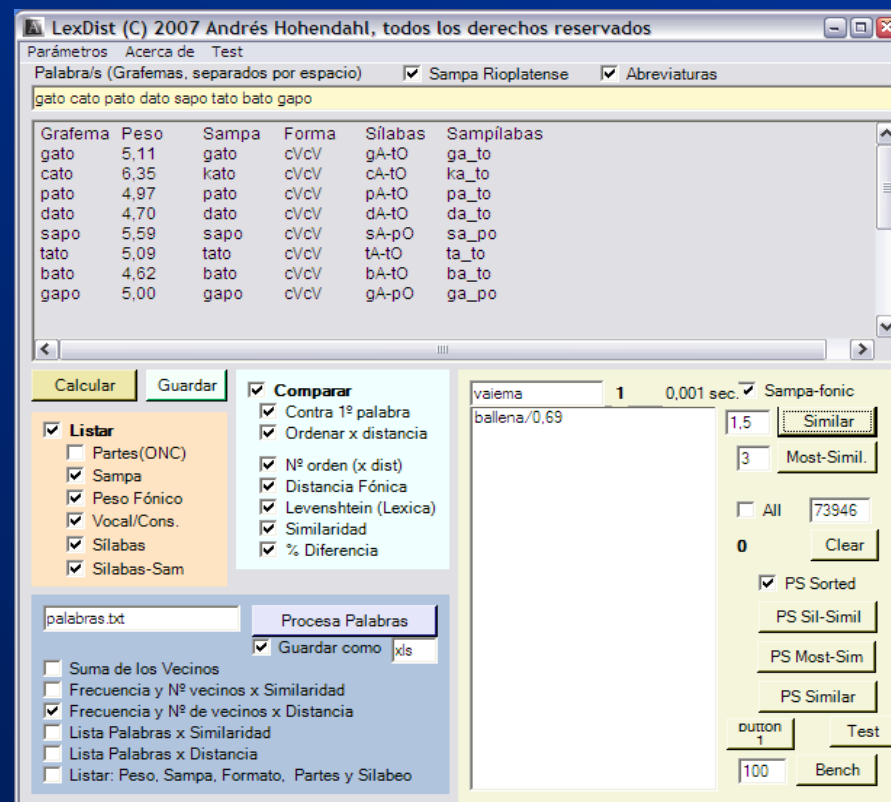
Benchmarking

    Fast-Find

    Best-Fit

Processing

    Word Lists

Making

    Word Matrix

# Features & Applications

## Features

- Find best-human like guess on mistyped or bad orthographically written (but sounding like) text.
- Delivering EAGLES 2.0 compatible, semantic-extended tags
- Uses Open-source Dictionaries, and spell checkers so it's adaptable to many languages based on free existing data.
- Delivers N-alternative Tags, ordered by phonetic distance.
- Detects foreign words  (tagging language) + capable to handle many mixed languages (one must be principal)

## Real World Human Computer Interface (HCI)

- Fast and Lightweight, engineered to fit into small appliances.
- Recognition + Guessing of parasynthetic O.O.V. (Out Of Vocabulary) in Scientific Text, Medical Records, etc.
- Robust Open-Lexicon Dialog System (free text)
- Automatic Speech Recognition (ASR) with huge-open Lexicon
- Teaching Aid / Support (Intelligent conversational agents)
- Artificial Understanding, AI, Context-based tagging, etc.

# Future Research Lines

H.C.I.

- Cognitive Modeling for fast Storage-Retrieval
- Spanish Dialog Subsystem
  - Robust GLR Compiler (Tomita-Like w/Scrödinger Tokens)
  - Cognitive Run-Time
    - Implied verbal Logic (Math, Set & Boolean Logic)
    - Simple Scientific Math (numeric + algebraic)
    - Scientific Units Cognitive Operations
    - Artificial Shallow Understanding
  - Information Extraction on OOV. & mistyped words (morphologically correctly constructed, even with errors)
  - Conversational-Space Resolution (Me-You-They)
  - On-The Fly Anaphora Resolution & used as context
  - Ontology Driven Contextual Conditional Parsing

Questions?

# Thank you

(dissertant)

José Francisco Zelasco

jfz@fi.uba.ar

Andrés T. Hohendahl

andres.hohendahl@fi.uba.ar