

## Capítulo 14

# Conversor de grafemas para fones baseado em regras para português

Sara Candeias e Fernando Perdigão

Esta apresentação tem por objectivo descrever um sistema de conversão automática de grafema para fone (GR2PH) para o português de Portugal. Para o desenvolvimento do GR2PH está a ser usado o corpus de unidades acentuais (palavras) em língua portuguesa SPEECHDAT (SPEECHDAT), disponibilizado pela Universidade do Minho (proveniente da colaboração entre a Linguateca e o Projecto Natura). A avaliação do GR2PH fará uso do vocabulário da base de dados SPEECHDAT bem como de outros corpora de teste já usados por diversos investigadores a trabalhar neste domínio. A anotação fonética de corpora em língua portuguesa seria um interessante recurso linguístico a tornar público na Linguateca. Este recurso poderia ficar disponível, depois de avaliado e validado o sistema.

A crescente procura de soluções baseadas em produtos de tecnologia da fala tem sido uma motivação para o desenvolvimento de sistemas capazes de estabelecer um interface Homem-Máquina mais natural, como são exemplos as práticas subjacentes a áreas do ensino/aprendizagem do português e da linguística clínica.

A consciencialização da necessidade destes produtos mobilizou ao desenvolvimento do GR2PH, que convertesse, de forma automatizada, corpora grafados em corpora notados foneticamente.

O GR2PH, do qual fazem parte os subsistemas ‘divisor de sílabas’ e ‘marcador de tonicidade’, é aquele para o qual o conhecimento linguístico contribui com um maior impacto. A estratégia adoptada para o GR2PH baseia-se em regras linguísticas cotejadas na estrutura da língua portuguesa. Para o desenvolvimento quer do sistema que transmuta grafema em fone, quer dos sistemas intermédios para divisão silábica e para marcação de sílaba tónica, foi usado o corpus de unidades acentuais (perto de 680000) em língua portuguesa, disponibilizado como recurso nascido da colaboração entre a Linguateca e o Projecto Natura. Na verdade, o acesso a este recurso resultou numa mais valia ao desempenho do(s) sistema(s) que se pretendia(m) desenvolver, e os testes que foram sendo feitos, mesmo de forma faseada, mostraram-se basilares na fase de estruturação da arquitectura do(s) próprio(s) sistema(s), complementares e final.

Para o português de Portugal, alguns transcritores de grafema para fone baseados em regras surgem descritos em Almeida e Simões (2001); Braga e Resende Jr (2007); Teixeira et al. (2006); Gouveia et al. (2000); Viana e Andrade (1985). Para a implementação das regras, em certos grupos, é reconhecida a importância da identificação da unidade silábica (Almeida e Simões, 2001; Braga e Resende Jr, 2007; Teixeira et al., 2006; Gouveia et al., 2000); noutras, é usada a informação da tonicidade da vogal (Almeida e Simões, 2001; Braga e Resende Jr, 2007; Viana e Andrade, 1985). A indispensabilidade de desenvolvermos um novo sistema de conversão GR2PH para o português de Portugal advém de factores como a escassa partilha dos algoritmos dos sistemas já implementados (dos quais poder-se-ia partir para um esforço de melhoramento do sistema) e dos resultados dos testes de desempenho provenientes de estudo comparativos. Este artigo apresenta uma tessitura alternativa de

| Convenções | Significado            |
|------------|------------------------|
| C          | consoante              |
| V          | vogal                  |
| .          | divisor de sílaba      |
| '          | marcador de tonicidade |
| #          | fronteira final de UA  |
|            | ou                     |

Tabela 14.1: Convenções usadas nas regras para implementação.

regras linguísticas a serem aplicadas no GR2PH para o português de Portugal, aliando a pertinência da informação linguística de regras de silabificação e de marcação de tonicidade. Resultando o sistema final da configuração de dois subsistemas perspectivados em regras inerentes à língua, o esforço do investimento tem por objectivo a viabilidade de um conversor capaz de uma eficácia que torne dispensável o recurso a dicionários de excepções. A arquitectura do GR2PH é resultado da complementaridade da aplicação do conhecimento linguístico e da ciência de engenharia, parceria esta que se traduz num diálogo necessário a uma execução que se pretende optimizada e eficaz.

## 14.1 Arquitectura do sistema de conversão GR2PH

O GR2PH recorre ao uso de sistemas intermédios, como o de separação da unidade acentual (UA, palavras) em sílabas e o de marcação de sílaba tónica (e conseqüente delimitação de sílaba(s) pré-tónica(s) e de sílaba(s) pós-tónica(s)). A vantagem desta abordagem explica-se pelo facto de ela permitir resolver a quase totalidade de casos de escolha fonética que não seria a acertada se resultasse apenas da inserção dos fones (nomeadamente vocálicos) considerados a partir de inventários fonéticos não diferenciados, isto é, não ponderados nem silabicamente nem atendendo à tonicidade em âmbito contextual de UA.

Todas as regras foram implementadas inicialmente em Matlab e foram testadas no vocabulário da base de dados SPEECHDAT (SPEECHDAT) e no corpus de unidades acentuais disponibilizado pela Linguateca/Projecto Natura.

Esta segunda parte apresenta as especificidades dos subsistemas de divisão silábica, de marcação de tonicidade e do transcritor, de forma a se ter uma visão global do sistema geral de conversão GR2PH. Na tabela 14.1 figuram as convenções usadas nas regras para implementação.

### 14.1.1 Subsistema de divisão silábica

A estrutura deste subsistema assenta a) num modelo de regras de divisão de base ortográfica, b) na consideração de vogal como núcleo de sílaba e c) na consideração de alguns dígrafos como grafema singular ('ch', 'ss', 'lh', 'gu'+ 'i'|'e', 'qu'+ 'i'|'e', etc.). O algoritmo

| Sequência | Exemplo       | Sequência | Exemplo      | Sequência | Exemplo   |
|-----------|---------------|-----------|--------------|-----------|-----------|
| CCVCC     | trans.cre.ver | CVCC      | subs.cre.ver | VC        | ac.tu.ar  |
| CCVVC     | grãos         | CVVC      | mães         | VV        | eu        |
| CVCCC     | tungs.té.ni.o | VCVC      | achar        | V         | á.gua     |
| CCCV      | stre.sse      | VVC       | aus.cul.tar  | CVV       | pai       |
| CCVC      | trás          | VCC       | abs.tra.ir   | CVC       | a.cam.par |
| CCVV      | grão          | CCV       | a.cre        | VC        | ac.tu.ar  |

Tabela 14.2: Lista dos padrões de sequências de grafemas a formar sílaba em português de Portugal.

do ‘divisor de sílabas’ reproduz uma busca feita por padrões de até 5 grafemas, resultando em 18 possíveis encontros de sequências que formam sílaba em português de Portugal (tabela 14.2). As regras foram distribuídas por dois grandes grupos para cada padrão de sequência de grafemas, isto é, considerando se na sílaba da UA a analisar é pertinente a informação dos 4 caracteres ou de mais que os 4 caracteres da sequência. Nesta repartição, surgem regras explícitas que apresentam um tipo repetido subsequente da iteração de sequências, como é exemplo a sequência VV presente nos padrões CCVV, CVVC, CVV e VVC. Na tabela 14.3, a título de exemplificação de procedimentos, surgem descritas regras para o padrão CVVC.

### 14.1.2 Subsistema de marcação de tonicidade

Na estruturação deste subsistema, toda a unidade (palavra) foi considerada acentual (UA) e, por isso, não foram admitidos segmentos desprovidos de tonicidade (Candeias, 2007). O algoritmo de marcação da sílaba tónica funciona com regras instituídas a partir da divisão silábica. Admitiu-se o acento tónico como o acento da UA (o acento principal), pelo que, nesta estrutura, não se considerou pertinente marcar os acentos secundários. Na tabela 14.4 figuram regras de marcação de sílaba tónica.

### 14.1.3 Subsistema de transcrição para fones

Para a anotação fonética, seguimos o alfabeto SAMPA para o português (SAMPA), sem o recurso a extensões como seria o caso das «oclusivas orais sonoras» «fricatizadas», traço que advém da posição em início de sílaba e intervocálica. Ainda que se tenha em vista a construção de um sistema de síntese futuro, o que leva a ter em conta, entre outros aspectos, a natureza particular de cada som em contexto de co-articulação e/ou de sandhi, o facto deste mapeamento da transmutação grafema–fone ir ser adicionado a um modelo acústico baseado em trifones, anula a necessidade de uma anotação fonética mais estreita. Com este mesmo princípio, não foram consideradas como «semiconsonânticas» ‘j’ e ‘w’ as unidades vocálicas grafadas ‘í’(ou ‘e’) e ‘u’ (ou ‘o’) dos ditongos ditos crescentes (pre-

|           | C           | V                | V      | C         | Grafema final da UA | Grupo silábico | Exemplo   |
|-----------|-------------|------------------|--------|-----------|---------------------|----------------|---|
| Sequência |             | a e o u<br>a e o | i<br>u | l r m s j | <b>V</b>            | <b>CVV.C</b>   | <i>pau.lada, mou.ro, tei.ma,<br/>lou.sa, bei.jo</i> |
|           |             | ã õ<br>ã         | e<br>o | ≠s #      |                     |                | <i>mãe.zinha, mão, ta.lão</i>                       |
|           | g q         | u                | a o    |           | <b>V</b>            |                | <i>quo.ciente, gua.rida,<br/>qua.se, qua.lidade</i> |
|           |             | a e o u<br>a e o | i<br>u | l z       | <b>#</b>            | <b>CV.VC</b>   | <i>pa.ul, ra.iz</i>                                 |
|           |             | a e o u<br>a e o | i<br>u | r m       | <b>C #</b>          |                | <i>ca.ir, ru.im, co.imbra</i>                       |
|           |             | a e o u          | i      | nh        | <b>V</b>            |                | <i>ba.inha, ta.bu.inha,<br/>mo.inho</i>             |
|           |             | a e o u<br>a e   | i<br>u | n         | <b>C</b>            |                | <i>re.incide, tran.se.unte</i>                      |
|           |             | a e o u<br>a e o | i<br>u | s         | <b>C #</b>          | <b>CVVC.</b>   | <i>ca.is, faus.to, a.zuis, bois</i>                 |
|           |             | ã õ<br>ã         | e<br>o | s         |                     |                | <i>mãos, pães</i>                                   |
|           | g q         | u                | a o    | l n r     | <b>C #</b>          |                | <i>qual, qual,quer, guar.da,<br/>quan.do</i>        |
|           | por defeito |                  |        |           |                     | <b>CV.VC</b>   | <i>be.ata, fi.os</i>                                |

Tabela 14.3: Ilustração de algoritmo de divisão silábica para o padrão de grafemas CVVC.

|               | Regra   | Marcador de tonicidade | Exemplo   |
|---------------|---|------------------------|---|
| <b>1.</b>     | Se na sílaba existirem vogais com acento gráfico  | sílaba em questão      | a.'ná.li.se                                       |
| <b>2.</b>     | Se na sílaba não existirem vogais sem acento gráfico  |                        |   |
| <b>2.1.</b>   | Se a UA tiver 1 sílaba  | sílaba em questão      | 'voz  |
| <b>2.2.</b>   | Se a UA tiver $\geq 2$ sílabas  |                        | pa.'ul<br>ra.'iz<br>ca.'ir                        |
| <b>2.2.1.</b> | Se for a última sílaba da UA com estrutura de<br>a e i o u + l r z<br>i u + $\emptyset$  s<br>i + m | sílaba em questão      | an.'dou, ca.pi.'tais<br>pe.'ru, pe.'rus<br>ru.'im |
| <b>2.2.2.</b> | por defeito   | penúltima sílaba       | a.na.'li.se                                       |

Tabela 14.4: Algoritmo de marcação de sílaba tónica.

| Fone | Posição de tonicidade | Posição silábica                      | Exemplos                               |
|------|-----------------------|---------------------------------------|--|
| o~   |                       | + m n (mesma sílaba)                  | 'om.bro → o~bru; pon.tu.'al → po~tual  |
| w~   |                       | ã + (mesma sílaba)                    | 'cão → k6~w~; cã.o.'zi.nho → k6~w~ziJu |
| o    | tónica                | + nh (sílabas seguintes)              | ri.'so.nho → rizoJu                    |
| O    | tónica                | + x (sílabas seguintes)               | pa.ra.'do.xo → p6r6dOksu               |
| o    | tónica                | + i (mesma sílaba)                    | 'oi.to → ojtú                          |
| o    | tónica                | + r (mesma sílaba e final de UA)      | pa.ssa.'dor → p6s6dor                  |
| O    | tónica                | + r (mesma sílaba)                    | 'cor.ta → kOrt6                        |
| o    | tónica                | + a (sílabas seguintes e final de UA) | 'to.da → tod6                          |
| O    | tónica                | por defeito                           | 'o.de → Od@; 'co.rre → kOR@            |
| O    | átona                 | (inicial de UA) + r                   | Or.ga.'ni.za → Org6niz6                |
| u    | átona                 | + r (mesma sílaba)                    | cor.'tar → kurtar                      |
| O    | átona                 | (inicial de UA)                       | o.'ní.ri.co → Oniriku                  |
| u    | átona                 | o (sílabas anteriores) +              | co.o. pe.ra.'ção → kuup@r6s6~w~        |
| u    | átona                 | (final de UA)                         | 'fi.lho → fiLu                         |
| O    | átona                 | + c p (mesma sílaba)                  | oc.'ta.vio → Otaviu; op.'ção → Ops6~w~ |
| u    | átona                 | por defeito                           | po.'ção → pus6~w~                      |

Tabela 14.5: Ilustração de algoritmo de conversão do grafema 'o' para fones.

sentes em relógio e em área, em suave e em nódoa). O algoritmo da conversão do grafema em fone funciona a partir das sílabas com 'marcação de tonicidade'. Isto é, a partir de um contexto-base, resultam casos de grafemas admitidos à conversão em fones que consideram a pertinência de informação da a) posição de tonicidade e da b) posição no âmbito da sílaba (na qual é pertinente o comportamento fonético dados os grafemas vizinhos). Na tabela 14.5 são exemplificados os algoritmos de conversão do grafema 'o' para os fones [o~], [w~], [o], [O] e [u], que resultam da atenção aos parâmetros descritos.

A análise e verificação de muitas regras foi conseguida por análise exaustiva ao corpus de UAs disponibilizado pela Universidade do Minho. Transcrições ou pronúncias alternativas não são consideradas neste sistema, como é o caso de homógrafos heterófonos.

## 14.2 Conclusão e trabalho futuro

Até esta fase, a forma gráfica convertida automatizadamente em forma fonética foi avaliada com referência à anotação manual. Dispomos apenas do vocabulário associado à base de dados SPEECHDAT como material de teste, embora a avaliação com este corpus não esteja ainda concluída, especialmente devida à discordância encontrada na conversão das semiconsoantes dos ditongos crescentes. Uma forma alternativa de fazer a avaliação do sistema consiste em comparar os resultados de vários sistemas de conversão – pelo menos um é de domínio público (Almeida e Simões, 2001) –, contando e analisando as diferenças encontradas. Como trabalho futuro, pretendemos construir uma aplicação on-line de conversão de grafemas para fones bem como de um corpus anotado foneticamente.