

**C-ORAL-BRASIL I**  
*Corpus* de referência  
do português brasileiro  
falado informal

UNIVERSIDADE FEDERAL DE MINAS GERAIS

REITOR Clélio Campolina Diniz

VICE-REITORA Rocksane de Carvalho Norton

EDITORA UFMG

DIRETOR Wander Melo Miranda

VICE-DIRETOR Roberto Alexandre do Carmo Said

CONSELHO EDITORIAL

Wander Melo Miranda (PRESIDENTE)

Antônio Luiz Pinho Ribeiro

Flavio de Lemos Carsalade

Heloisa Maria Murgel Starling

Márcio Gomes Soares

Maria das Graças Santa Bárbara

Maria Helena Damasceno e Silva Megale

Roberto Alexandre do Carmo Said

*Tommaso Raso*

*Heliana Mello*

Organizadores

# C-ORAL-BRASIL I

*Corpus* de referência  
do português brasileiro  
falado informal

Belo Horizonte  
Editora UFMG  
2012

© 2012, Os autores  
© 2012, Editora UFMG

Este livro ou parte dele não pode ser reproduzido por qualquer meio sem autorização escrita do Editor.

---

C787 C-ORAL-BRASIL I : *corpus* de referência do português brasileiro falado informal / Tommaso Raso, Heliana Mello, organizadores. – Belo Horizonte : Editora UFMG, 2012.  
332 p. : il + 1 DVD-ROM. – (Linguajar)

Acompanha DVD com conteúdo completo do livro, arquivos de alinhamento, anotação gramatical, especificações do *corpus*, medidas estatísticas, listas de frequência e arquivos de som.

Inclui bibliografia.

ISBN: 978-85-7041-943-9

1. Linguística. 2. Comunicação oral. I. Raso, Tommaso.  
II. Mello, Heliana. III. Série.

CDD: 410

CDU: 81

---

Elaborada pela DITTI – Setor de Tratamento da Informação  
Biblioteca Universitária da UFMG

Este livro recebeu apoio financeiro da Fapemig.

COORDENAÇÃO EDITORIAL Danívia Wolff

ASSISTÊNCIA EDITORIAL Eliane Sousa e Euclídia Macedo

COORDENAÇÃO DE TEXTOS Maria do Carmo Leite Ribeiro

PREPARAÇÃO DE TEXTOS Alexandre Vasconcelos, Danívia Wolff e

Maria do Rosário Alves Pereira

REVISÃO DE PROVAS Beatriz Trindade, Camila Figueiredo, Davi Bezerra de Souza

COORDENAÇÃO, PROJETO GRÁFICO E FORMATAÇÃO Cássio Ribeiro

CAPA E PRODUÇÃO GRÁFICA Diêgo Oliveira

EDITORA UFMG

Av. Antônio Carlos, 6.627 – Ala direita da Biblioteca Central – térreo

Campus Pampulha – 31270-901 – Belo Horizonte / MG

Tel.: + 55 31 3409-4650 Fax: + 55 31 3409-4768

www.editora.ufmg.br editora@ufmg.br

Para o Nando  
(*in memoriam*).



---

## Sumário

Prefácio	Emanuela Cresti   Massimo Moneglia	13
Apresentação	Tommaso Raso   Heliana Mello	27
CAPÍTULO 1	Heliana Mello	
Os <i>corpora</i> orais e o C-ORAL-BRASIL		31
1. Introdução		31
1.1 Os <i>corpora</i> orais no cenário internacional		34
1.2 Parâmetros para a compilação de <i>corpora</i> orais		36
2. Os bancos de dados e <i>corpora</i> orais brasileiros		38
2.1 NURC - Projeto Norma Linguística Urbana Culta		41
2.1.1 NURC-SP		43
2.1.2 NURC-RJ		43
2.1.3 NURC-Salvador		44
2.1.4 NURC-Recife		44
2.1.5 NURC-Porto Alegre		45
2.2 Projeto Discurso e Gramática		45
2.3 VARPORT: Análise Contrastiva de Variedades do Português		45
2.4 PROFALA: Variação e Processamento da Fala e do Discurso: análises e aplicações		47
2.5 VALPB: Projeto Variação Linguística no Estado da Paraíba		47
2.6 Projeto Vertentes do Português Popular do Estado da Bahia		48
2.7 VARSUL: Variação Linguística na Região Sul do Brasil		48
2.8 ALIP: Projeto Amostra Linguística do Interior Paulista – o banco de dados IBORUNA		49
3. O C-ORAL-BRASIL		50
4. Considerações finais		53

O <i>corpus</i> C-ORAL-BRASIL	55
1. Introdução	55
2. A arquitetura	58
2.1 A perspectiva pragmática: a variação diafásica	61
2.2 A variação diastrática	64
2.3 Outros aspectos da arquitetura	70
2.3.1 A variação diatópica	70
2.3.2 Considerações sobre alguns arquivos de gravação	71
2.3.3 Um <i>minicorpus</i> etiquetado informacionalmente	73
3. Questões metodológicas	73
3.1 A qualidade acústica	73
3.2 As segmentações	75
3.3 As transcrições	79
4. Uma comparação com o C-ORAL-ROM	80
4.1 A arquitetura	80
4.2 A metodologia	86
4.2.1 As gravações	86
4.2.2 A transcrição e a segmentação	87
4.2.3 O <i>parser</i>	89

O C-ORAL-BRASIL e a Teoria da Língua em Ato	91
1. Introdução	91
2. Enunciados e atos de fala	92
2.1 O enunciado	92
2.2 Os atos de fala	99
3. A unidade tonal e a unidade informacional	105
3.1 As unidades textuais	107
3.2 As unidades dialógicas	110
4. O limite do isomorfismo	113
4.1 As unidades de escansão	113
4.2 Os comentários múltiplos	114
4.3 As estrofes	117
5. Anotações de sintaxe	121



Transcrição e segmentação prosódica do *corpus* C-ORAL-BRASIL:

critérios de implementação e validação	125
1. Visão geral	125
2. A formação dos transcritores	126
2.1 Os grupos de transcritores	127
2.2 Treinamento para a segmentação prosódica	128
2.3 Avaliação do processo de constituição do time de transcritores	130
3. A transcrição	130
3.1 Sobre os critérios adotados na transcrição	133
3.2 Convenções gerais e formas transcritas ortograficamente	134
3.2.1 Ruídos paralinguísticos	134
3.2.2 Hesitações e palavras interrompidas	135
3.2.3 Onomatopeias	136
3.2.4 Interjeições e exclamações	136
3.2.5 Siglas e acrônimos	137
3.2.6 Numerais	138
3.2.7 Palavras estrangeiras e erros de pronúncia	138
3.2.8 Formas transcritas conforme a ortografia padrão	139
3.2.9 Palavras não transcritas ou censuradas	139
3.3 Critérios não ortográficos	140
3.3.1 Aférese	140
3.3.2 Fenômenos relativos à conjugação verbal	140
3.3.3 O plural	141
3.3.4 Fenômenos relativos a pronomes	142
3.3.5 Preposições articuladas e reduzidas	142
3.3.6 Articulação de preposições e pronomes	142
3.3.7 Negação	143
3.3.8 Construções interrogativas, o pronome relativo e pseudorrelativo	143
3.3.9 As formas <i>senhor</i> e <i>senhora</i>	144
3.3.10 Diminutivos	144
3.3.11 O intensificador <i>mó</i>	144
3.3.12 Rotacismo	145
3.4 Algumas observações sobre a transcrição da fala espontânea	145
4. A segmentação prosódica	146
4.1 Convenções da segmentação prosódica	147
4.2 Procedimentos adotados na anotação da segmentação	148

5.	Validação de <i>corpora</i> orais e o caso do C-ORAL-BRASIL	150
5.1	A validação da segmentação prosódica	152
5.1.1	Realização de testes	153
5.1.2	Resultados da validação da segmentação: grupo 1	156
5.1.3	Resultados da validação da segmentação: grupo 2	160
5.1.4	Reavaliação do acordo entre os transcritores do grupo 1	164
5.1.5	Considerações gerais sobre a validação da segmentação prosódica	166
5.2	A validação da transcrição	167
5.2.1	Método	169
5.2.2	Resultado das validações prévia e final	171
6.	Observações finais	176

CAPÍTULO 5 | Tommaso Raso | Maryualê M. Mittmann

	As principais medidas da fala	177
1.	Objetivos do capítulo	177
2.	Uso de verbos e nomes na fala espontânea e na escrita	178
2.1	Ocorrência de nomes e verbos na fala e na escrita do PB	180
2.2	Análise interlinguística da ocorrência de nomes e verbos na fala espontânea	182
3.	Unidades estruturais naturais da fala espontânea	185
3.1	O turno dialógico	185
3.2	Estruturação da fala em enunciados simples e complexos	188
3.3	Complexidade informacional na fala espontânea	192
3.4	Fenômenos de fragmentação da fala	195
3.5	Tamanho e duração das unidades de referência da fala	198
4.	Enunciados verbais e não verbais	201
5.	Conjunções coordenativas e subordinativas na fala e na escrita	206
5.1	A conjunção <i>e</i>	212
5.2	A conjunção <i>mas</i>	214
5.3	A conjunção <i>que</i>	215
5.4	A conjunção <i>porque</i>	217
6.	Considerações finais	219
	Anexo	220

A anotação gramatical do C-ORAL-BRASIL	223
1. Introdução	223
2. O ponto de partida: o <i>parser</i> PALAVRAS	225
3. A normalização do fluxo textual	230
4. A normalização do léxico e da ortografia	236
5. A segmentação sintática	243
6. A avaliação	247
7. Conclusão	249
8. Apêndice: definições das etiquetas	250
Bibliografia	255
Anexos	269
Anexo 1 - Lista de frequência das primeiras 100 formas no <i>corpus</i> C-ORAL-BRASIL e no Banco do Português	269
Anexo 2 - Lista de frequência dos primeiros 100 lemas no <i>corpus</i> C-ORAL-BRASIL	273
Anexo 3 - Lista de frequência dos primeiros 100 verbos no <i>corpus</i> C-ORAL-BRASIL, em comparação com o <i>corpus</i> do Português Europeu do C-ORAL-ROM	275
Anexo 4 - Lista de frequência dos primeiros 100 nomes no <i>corpus</i> C-ORAL-BRASIL, em comparação com o <i>corpus</i> do Português Europeu do C-ORAL-ROM	278
Anexo 5 - Lista de frequência dos primeiros 100 adjetivos no <i>corpus</i> C-ORAL-BRASIL, em comparação com o <i>corpus</i> do Português Europeu do C-ORAL-ROM	282
Anexo 6 - Lista de frequência dos primeiros 100 advérbios no <i>corpus</i> C-ORAL-BRASIL, em comparação com o <i>corpus</i> do Português Europeu do C-ORAL-ROM	285
Anexo 7 - Lista de frequência das primeiras 10 conjunções e das primeiras 10 preposições no <i>corpus</i> C-ORAL-BRASIL	289
Anexo 8 - Formas não ortográficas	289
Sobre os autores	331



---

## Prefácio

É com grande satisfação que saudamos a publicação do *corpus* C-ORAL-BRASIL, fruto do trabalho intenso de Tommaso Raso e Heliana Mello e do entusiasmo do grupo de jovens do Laboratório de Estudos Empíricos e Experimentais da Linguagem (LEEL) da UFMG (Maryualê Mittmann, Heloísa Vale, Priscila Côrtes, entre outros).

Esta obra, que vimos começar em 2007-2008 e cujo desenvolvimento acompanhamos de perto dentro do acordo de colaboração entre a Università di Firenze e a UFMG, amplia, mesmo que de maneira independente, a coleção de línguas nacionais românicas da Europa (italiano, francês, português europeu, espanhol), representadas no C-ORAL-ROM, acrescentando a ela o português brasileiro (PB). Trata-se de um fato relevante para a fala românica em geral porque o C-ORAL-BRASIL permite uma representação mais definida do PB, à qual a comunidade científica poderá recorrer para marcar a distinção entre PB e português europeu (PE) – após séculos de contato do PB com línguas tão diferentes, em primeiro lugar ameríndias e depois africanas (cf. Mello, 1997 e 2011). É relevante também pela proeminência atual do Brasil na cena político-econômica internacional e pelo peso de 180 milhões de falantes, cujo uso linguístico pode ser retratado nacional e internacionalmente. Trata-se, enfim, de um fato relevante pelo nível de representação da fala alcançado nesta obra.

A escolha de desenvolver um novo recurso da fala espontânea brasileira seguindo o modelo do C-ORAL-ROM, em particular depois dos importantes trabalhos realizados no Brasil nesse domínio a partir dos anos de 1970 (como o NURC, entre outros), não era óbvia. Gostaríamos de lembrar aqui, rapidamente e do nosso ponto de vista, as motivações e o quadro cultural de referência que a bela descrição de Raso e o *background* ilustrado por Mello mostrarão em detalhe ao leitor mais à frente.

O *corpus* de fala românica C-ORAL-ROM é fruto de um projeto financiado pelo Quinto Programa Quadro da União Europeia, realizado por meio de um consórcio de entes públicos e privados (Università di Firenze, Universidad Autónoma de Madrid, Université Aix-Marseille en Provence, Centro de Linguística da Universidade de Lisboa, Fundação Bruno Kessler, Pitch France) concluído em 2004 com a publicação de um recurso multimídia que é distribuído pela indústria das línguas da European Language Resource Distribution Agency e na academia pela John Benjamins Publishing Company. O recurso permite o acesso simultâneo ao som e à transcrição de sessões gravadas das quatro coleções comparáveis de fala espontânea (italiano, espanhol, francês e português) e é acompanhado de instrumentos informáticos que permitem a análise acústica e a recuperação da informação textual para fins de estudo, aprendizagem e utilização para o desenvolvimento das tecnologias da linguagem.

A idealização do C-ORAL-ROM, no final dos anos de 1990, seguia uma série de relações científicas entre linguistas, em especial franceses, italianos e portugueses (junto conosco, Claire Blanche-Benveniste e Fernanda Bacelar do Nascimento, e, mais tarde, Francisco Marcos Marin), que reconheciam a urgência do estudo da realização falada da linguagem, a qual havia se tornado possível somente em tempos relativamente recentes, em função do desenvolvimento tecnológico. Esses estudiosos reconheceram, principalmente, o interesse teórico que o estudo apresentava, tanto para a descrição das línguas nacionais, quanto para os nossos conhecimentos gerais sobre a linguagem, tão diferente na sua realização oral em relação à sua idealização escrita, na qual sempre se basearam as gramáticas, e, a partir da segunda metade do século passado, a linguística da competência.

Essa exigência havia levado, nos 20 anos anteriores, à formação de vários recursos de fala, na Itália, em Portugal, na França e na Espanha, ligados a abordagens e tradições em parte diferentes, mas acomodadas pela perspectiva do estudo linguístico baseado em *corpora*. O C-ORAL-ROM deriva da necessidade de se colocarem à disposição dados mais estritamente comparáveis para o estudo da realização oral da linguagem em nível interlinguístico, sem os quais seria impossível evidenciar, com adequada base de observação, tanto as propriedades diferenciadoras da oralidade quanto os aspectos que, nessa forma de realização, diferenciam entre si as línguas naturais.

O C-ORAL-ROM fixou dois pontos nessa tarefa, definindo, por um lado, os requisitos nas modalidades necessárias para a representação do dado linguístico e, por outro, requisitos relativos à formação dos *corpora* para fins da sua comparabilidade. Ambos foram seguidos e tornados ainda mais eficazes na realização do C-ORAL-BRASIL, fazendo desse *corpus* um recurso plenamente compatível com seu antecedente. Vale a pena, aqui, ressaltar essas escolhas e a maneira exemplar como foram levadas adiante neste trabalho. Falaremos brevemente dos critérios de formação do *corpus* e com mais detalhe dos critérios de representação da fala. Ambos constituem a premissa para o desenvolvimento de estudos comparativos sobre a oralidade.

O requisito da comparabilidade na formação do *corpus* e, ao mesmo tempo, o requisito de representação da fala espontânea, juntos, constituem um desafio. Não é possível, de fato, constituir *corpora* paralelos e, portanto, forçadamente comparáveis, sem renunciar à característica principal da fala – a espontaneidade. A fala espontânea é caracterizada por grande variação, em conexão seja com as características do falante, seja com as características da situação em que acontece o ato comunicativo. Em particular, a fala se produz com base na interação, e a variedade dos atos de fala realizados na vida cotidiana e das suas situações eliciantes é ampla e desconhecida.

O C-ORAL-ROM respondeu a esses requisitos seguindo a regra de ouro da Linguística de *Corpus* para a formação de *corpora* de referência (veja-se o British National Corpus - BNC): definição dos critérios de variação utilizados para a representação do universo e assunção dos mesmos critérios na constituição de cada *corpus*.

As linhas ao longo das quais a fala varia são, contudo, muitas. A escolha dos critérios de amostragem é, portanto, um capítulo essencial para assegurar a representatividade dos recursos, seja em nível intralinguístico (representar a variação), seja em nível interlinguístico (representá-la da mesma maneira em cada língua).

Já o BNC, na sua seção de oralidade, individualizava dois critérios de amostragem a serem seguidos em paralelo: amostragem por contexto e amostragem por falante. De um lado, coletam-se amostras de fala de segmentos diastraticamente balanceados da população; de outro, identificam-se *a priori* contextos típicos nos quais a qualidade da produção oral é claramente diferente (bate-papos, fala oficial ou acadêmica, fala midiática,

entrevistas). Do cruzamento dos dois parâmetros, espera-se conseguir uma representação suficiente do universo.

Infelizmente, a prática mostrou que essa escolha, na realidade, reduz aos simples bate-papos o âmbito em que a variação da fala cotidiana é maior, ou seja, o uso informal da linguagem. A amostragem diastrática, de fato, produz amostras unicamente das situações mais fáceis de serem gravadas: os bate-papos em situação privada e familiar. Como a dificuldade do *corpus* israelense, o CoSIH (Izre'el; Hary; Rahav, 2001), mostrou, não se consegue, por muitos motivos que lembramos em várias ocasiões, gravar os falantes em todas as situações por eles encontradas durante o dia. Isto reduz enormemente a possibilidade de documentar, no *corpus*, a variedade dos contextos que, na vida cotidiana, determinam exigências interativas diferentes e produzem a exigência de atividades linguísticas as mais variadas. É essa propriamente a especificidade da fala espontânea.

Em outras palavras, na prática, a amostragem diastrática não produz variação diafásica, e, contrariamente a qualquer critério de amostragem, a variação diafásica nos contextos informais tem probabilidade zero de ocorrência, com a exceção dos bate-papos.

A escolha do C-ORAL-ROM e, com sucesso bem maior, a escolha do C-ORAL-BRASIL tentaram modificar o estado da arte dando espaço à maior variação diafásica possível. Por um lado, renunciou-se a um balanceamento diastrático da amostra, por outro, escolheram-se fixar critérios gerais que aumentassem as probabilidades de ocorrência do maior número possível de situações de interação específicas do uso informal da linguagem: a) variação por estrutura do evento comunicativo (diálogo, monólogo, multidialogo); b) variação por contexto social em que o evento ocorre (vida pública *versus* vida privada e familiar). Cruzando os dois critérios e requerendo a maior variação possível de situações em cada ramificação da amostragem, conseguimos documentar o informal muito além dos bate-papos, com muitas amostras de fala fortemente interativa em situação (enquanto se conserta um carro, enquanto se briga, enquanto se dá aula de direção, enquanto se revelam fotografias, enquanto se tenta entrar em acordo com os amigos sobre o que comprar e como pagar, enquanto se lida com clientes em um bar, no guichê dos impostos, no trabalho etc.).

A documentação desse tipo de variação constitui o verdadeiro valor adjunto dos *corpora* de fala espontânea e, no estado atual, pode ser documentada em um *corpus*, e com muito trabalho, somente com base em



uma escolha precisa de amostragem. Os estudantes e os professores que os orientam, e é esse o caso do C-ORAL-BRASIL, são os vetores privilegiados que descobrem novas possibilidades de documentação, penetrando em seus ambientes de origem, necessariamente os mais variados em sociedades complexas. A casualidade requerida pelas amostragens desse tipo é limitada, portanto, pelas restrições do método, mas não tem como consequência reduzir a zero a probabilidade de ocorrência das situações de interação em que se usa a oralidade informal além dos bate-papos.

Acreditamos que o C-ORAL-BRASIL alcançou o maior resultado possível com essa metodologia, ultrapassando em qualidade e quantidade a variação diafásica dos *corpora* documentados no C-ORAL-ROM. Menos de 20% dos textos de interação dialógica ou conversacional no C-ORAL-BRASIL são constituídos por bate-papos ou entrevistas, enquanto grande parte das sessões de gravação documentam situações caracterizadas por acionalidade em movimento (veja-se o Capítulo 2), tornando possível até a gravação de um grupo de amigos enquanto jogam uma partida de futebol.

Dos muitos problemas que podem surgir em uma amostragem desse tipo, o C-ORAL-BRASIL evitou provavelmente o maior: documentar falantes pertencentes à mesma camada social, achatando a diversidade linguística que é função da diastratia.

Como se evidencia no Capítulo 2, o balanceamento diastrático, resultado do amplo número de falantes registrados no C-ORAL-BRASIL, mesmo que não perseguido intencionalmente, é um notabilíssimo íon. Ao contrário de como normalmente se pensa, na prática, não é tanto a variação diastrática que induz a variação diafásica, mas é esta última que induz a primeira.

Ressaltamos, então, não sem uma ponta de inveja, que o C-ORAL-BRASIL constitui um passo significativo no estado da arte da constituição dos *corpora* orais do ponto de vista dos resultados obtidos na amostragem do universo. Será preciso muito tempo para que a comunidade de estudiosos considere as consequências da disponibilidade desses dados para o estudo da linguagem e se adéque a esse “passo à frente”.

O primeiro aspecto é de fácil compreensão. Não é possível estudar a oralidade perdendo a sua característica informação acústica ou sem uma representação capaz de fixar a volatilidade do evento fonético em unidades informacionais mínimas, as palavras. A solução não é, contudo, de realização fácil, porquanto necessita, por um lado, da adoção de formatos textuais

computáveis adequados à representação do diálogo e, por outro lado, de critérios de relação entre o evento fonético e a sua contraparte linguística. Não por acaso, a maior parte dos *corpora* de fala realizados nos anos de 1990 não permitiam o acesso ao sinal acústico, e mesmo a qualidade das transcrições era deficitária se comparada com a sua fonte.

Precisamos lembrar que a concepção dos *corpora* de fala mudou totalmente durante os últimos 20 anos. De fato, ainda no começo dos anos de 1990, a sua relevância era concebida essencialmente como possibilidade de acesso a recursos lexicais e, só parcialmente, sintáticos diferentes da escrita. Disto constituem prova os grandes *corpora* da época, dentre os quais o primeiro é o BNC, cuja estruturação se mantém ainda hoje exemplar, mas que, em sua seção de fala, que previa 10 milhões de *tokens* sobre os 100 do recurso total, não previa o acesso à contraparte acústica (disponível agora depois de duas décadas), e cujas transcrições, é um fato notório, foram feitas praticamente por funcionários e não por pesquisadores treinados. O mesmo pode ser dito do *corpus* alemão (Mannheim Corpus), dos anos de 1970, e, pelo que sabemos, do grande *corpus* brasileiro NURC. Todos estes apresentam, portanto, ausência total ou parcial do correspondente sonoro da gravação, que, se feito em suporte analógico, pode se desgastar completamente em uma década, além de transcrições ortográficas frequentemente parciais, realizadas com base em critérios não declarados, às vezes resumindo, às vezes normalizando, e raramente feitas com acurácia e em conformidade com o dado sonoro originário.

No C-ORAL-BRASIL, assim como no C-ORAL-ROM, a transcrição foi realizada segundo critérios ortográficos, ou seja, não fonéticos, e seguindo formatos padrão já a partir dos anos de 1980 (CHAT). A exploração de práticas consolidadas como o formato CHAT evita a introdução de convenções idiossincráticas na representação textual, que reduzem a possibilidade de utilização do recurso na comunidade científica e maximizam as possibilidades de computação automática da informação textual.

Vale então a pena lembrar, brevemente, algumas implementações desse formato que se confirmaram funcionais para a constituição de grandes *corpora* de fala românica. Primeiramente, a manutenção de uma sequência unitária para preservar o turno dialógico de cada falante, que deve ser entendido como a primeira entidade natural da fala; e em segundo lugar, a anotação sistemática, dentro da transcrição em nível zero, ou seja, no nível da transcrição da palavra, do que foi definido como *quebras prosódicas terminais*

e *não terminais* (//, /). As quebras indicam todas as rupturas perceptualmente relevantes para a avaliação da organização das produções faladas em agrupamentos finalizados do ponto de vista informacional.

A escolha da representação ortográfica, e não fonética, implica o desejo de estudar a oralidade com base em seu valor linguístico e não com base em sua primeira articulação. Nisto, o C-ORAL-BRASIL se encontrou em um ponto de fronteira difícil, e o valor e o interesse do que foi proposto devem ser ressaltados. Se a transcrição fonética torna obscuros os elementos da cadeia da fala aos quais atribui valor linguístico, seguir a ortografia padrão definida com base em uma idealização abstrata da língua tornaria inutilizável o dado vivo da fala para evidenciar os processos morfossintáticos próprios do português brasileiro, obscurecendo processos dinâmicos que caracterizam componentes inteiros dessa variedade, a morfologia livre e flexiva em primeiro lugar.

A transcrição da fala de línguas como o italiano, o francês, o espanhol e o PE pôde se servir de convenções de transcrição, mesmo se não compartilhadas unanimemente, de variedades regionais afirmadas e consolidadas e, de toda maneira, não muito distantes do padrão nacional. Não era esse, contudo, o caso do *corpus* de fala brasileiro. A questão envolve vários níveis: em primeiro lugar, a coleção brasileira é de fala mineira, mais especificamente, da cidade de Belo Horizonte e de sua área metropolitana – tal como requisitado por cada *corpus* de fala, que deve ser representativo de uma única variedade regional. Parece que muitas tendências gerais em curso na fala brasileira foram primeiro atestadas exatamente nessa variedade regional. Portanto, nos encontraríamos frente a uma situação de “vanguarda” linguística do mineiro que, evidentemente, deve ser registrada. Entretanto, não é fácil distinguir entre o que é estável dentro de uma variedade e, portanto, relativo ao sistema, e o que pode ser uma variante idiossincrática de um falante, que, ao contrário, deve ser normalizada, para não comprometer a reconhecibilidade.

Entretanto, o problema mais importante é que a fala mineira, como uma das variedades mais significativas do PB, é muito distante da de Lisboa ou *tout court* do PE. Não compete a nós, neste Prefácio, enfrentar a questão de qual é a relação entre PB e PE. Certamente, tivemos a possibilidade de verificar pessoalmente a distância entre um e outro. Uma transcrição da fala mineira, por mais fiel que seja ao conteúdo linguístico de entidades terminadas alinhadas com o som, impõe o problema de o quê e como transliterar.

O sistema pronominal sujeito nas suas variantes tônica e átona, a falta da marca de plural em nomes e adjetivos, a redução da morfologia verbal, a nova declinação verbal para alguns lemas ou tempos, o léxico diferente, mas também palavras existentes em PE cuja forma não é mais a mesma, fórmulas, interjeições, alocutivos autóctones, ou seja, a condição de ter que refletir o som por um lado e o significado por outro induziu os autores a se confrontarem com um novo sistema complexo de escrita ortográfica, tarefa de peso significativo. Sentimos vontade de dizer, a esse respeito, que as soluções escolhidas nos parecem muito conscientes e compartilháveis, além de testemunharem coragem e cautela ao mesmo tempo. O resultado é uma transcrição correspondente ao som originário, legível e compreensível para qualquer falante brasileiro, que, contudo, achará traços específicos da variedade mineira que podem não ser os seus. Provavelmente a transcrição parecerá também escandalosa a quem se atém às convenções do português escrito padrão, e podemos imaginar que não faltará quem diga que os textos assim transcritos são obscuros e incompreensíveis. Já aconteceu, para os textos italianos falados na variedade toscana, que a presença de algum apóstrofo para sinalizar aférese e truncamento, algumas poucas variantes na morfologia verbal, algum lema regional, locução verbal ou fraseologia não pan-italiana, mas absolutamente transparente, levasse o *corpus* a ser considerado como dialetal e sem validade representativa.

As convenções, nos casos em que se decidiu por uma representação ortográfica não padrão da palavra no C-ORAL-BRASIL, ilustradas no Capítulo 4, foram fruto de reflexão intensa por parte dos organizadores (veja-se também Mello e Raso, 2009). Elas são, ao mesmo tempo, garantia da qualidade das transcrições e uma premissa para a utilização dos dados do *corpus* oral para uma descrição gramatical. Voltaremos a esse último ponto, mas, quanto às garantias sobre a confiabilidade da transcrição, o C-ORAL-BRASIL apresenta um resultado absolutamente original no domínio dos *corpora* orais. A adequação da transcrição foi testada e é garantida não somente, como normalmente é feito, com relação à sua computabilidade (vejam-se os resultados relativos à etiquetagem morfossintática no Capítulo 6), mas também frente ao dado acústico. Os resultados dos testes de validação, tanto gerais quanto nas formas não padrão, evidenciam um nível de adequação além de qualquer requisito: uma margem de erro geral inferior a 0,81% e inferior a 0,57% para as formas relativas aos critérios não padrão implementados pelos autores,

sendo que o fenômeno que apresentou a maior porcentagem de erro (as preposições apostrofadas) não superou os 3,28% de erros (veja-se o Capítulo 4). Claire Blanche-Benveniste teria gostado desse nível de acurácia, ela que foi a primeira a transferir uma tradição filológica de precisão e transparência nos critérios de transliteração. Os linguistas que desejarem utilizar o *corpus* C-ORAL-BRASIL para identificar as mais consistentes mudanças do PB com relação ao PE poderão, portanto, confiar nos dados validados quanto à uniformidade ortográfica dos traços morfológicos característicos.

A nossa avaliação é que esse trabalho marca um limite a partir do qual existirá um *ante quem* e um *post quem*, a partir do qual se tornará mais fácil realizar também novas coleções de diferentes variedades brasileiras e conduzir pesquisas gramaticais que deverão levar em conta os dados linguísticos, que terão rapidamente mudado em uma sociedade em expansão como a brasileira. A impressão é de que, devido à influência que contatos com línguas, principalmente africanas, geraram no português e seu aprendizado como segunda língua e não como língua materna pela maior parte da população, o atual uso pervasivo do português deva ser associado a uma inevitável e profunda transformação, de que o *corpus* de fala mineira constitui um testemunho precioso.

Apesar de a escolha da transcrição ortográfica denunciar o interesse desse *corpus* pelo estudo linguístico, os foneticistas e os tecnólogos que queiram utilizá-lo para o estudo da voz não ficarão decepcionados e poderão usufruir da altíssima qualidade da fonte acústica. A dedicação e a acurácia com as quais foram efetuadas as gravações possibilitaram um resultado até pouco tempo inimaginável para um grande *corpus* de fala espontânea, que é obrigado a gravar uma ampla variedade de situações diafásicas e contextuais, cada uma necessitando de soluções específicas de microfonação. O uso de microfones monodirecionais, praticamente um para cada falante, eliminou em grande parte os fenômenos de retorno em todos os textos, sempre identificados com um metadado relativo à qualidade acústica. Tivemos a possibilidade de ter acesso ao dado acústico de várias sessões de gravação incluídas no C-ORAL-BRASIL, e os espectrogramas quase não se diferenciam dos traçados obtidos com base em gravações em câmara anecoica. Por isso, o recurso resulta adequado à transcrição fonética e ao alinhamento silábico para os estudiosos que desejarem utilizá-lo com esse objetivo no futuro.

Na representação da fala, o alinhamento do som à transcrição é o ponto, teoricamente, mais delicado, sendo realizado no C-ORAL-ROM e no C-ORAL-BRASIL segundo um critério e uma técnica comuns.

O desenvolvimento vertiginoso das possibilidades técnicas permitiu transformar essa exigência naquilo que hoje nos parece um requisito correto para um bom usufruto e uma boa apreciação do texto oral. A técnica utilizada é fornecida pelo *software* WinPitch, de Philippe Martin, adaptado à função de alinhamento de grandes *corpora* para o projeto C-ORAL-ROM. O WinPitch permite o alinhamento em níveis independentes de porções do sinal e oferece a função de sincronia som-transcrição, acompanhada pela análise acústica do sinal em tempo real.

O alinhamento produz dois efeitos fundamentais e que, em parte, poderiam parecer contraditórios. Por um lado, a presença do som obriga a uma transcrição em conformidade com a modalidade original. Por outro, o fato de que uma transcrição de fala espontânea quase nunca pode ser considerada definitiva, porque a escuta e a reescuta induzem, quase sempre, a pequenas correções, leva à consciência de que ela não é inteiramente autossuficiente e deve sempre ser acompanhada pelo som. A qualidade da transcrição no C-ORAL-BRASIL depende, também, das revisões realizadas após o alinhamento.

O alinhamento implica também o problema de o quê é alinhado ao quê. A escolha do enunciado como unidade de alinhamento corresponde à escolha desse nível pragmático de descrição linguística como unidade de referência da língua falada (Cresti, 2000a e 2000b). Como bem argumenta Raso, no Capítulo 3, a Teoria da Língua em Ato oferece uma solução operacional para a identificação do enunciado, através da identificação das unidades concluídas por quebras prosódicas terminais (enunciados, padrões ilocucionários, estrofes). Estas constituem as unidades de referência alinhadas ao som: de silêncio ou quebra prosódica terminal ao sucessivo silêncio ou quebra prosódica terminal.

O critério de definição dos enunciados no C-ORAL-BRASIL, como no C-ORAL-ROM, é, portanto, prosódico e se fundamenta na percepção das rupturas entonacionais da cadeia falada percebidas como terminais pelo transcritores.

A introdução, na representação falada, dos critérios perceptuais, especificamente em pontos cruciais como a identificação das unidades de referência, essenciais para qualquer análise linguística superior ao nível lexical, frequentemente deixava perplexos muitos colegas, mesmo os

mais prestigiosos, que, por um lado, temiam a introdução de fatores idiossincráticos e não verificáveis e, por outro, estavam acostumados a utilizar para esses fins índices sintáticos (Biber) ou interpretativos (Blanche-Benveniste) sentidos como “mais linguísticos”.

Notar a correlação estrita entre unidades prosódicas e atividades pragmáticas é o ponto essencial que caracteriza a série de recursos C-ORAL-BRASIL / C-ORAL-ROM e corresponde à concepção mais geral da fala espontânea como interação pragmática, mais que como atividade de construção textual.

A correlação entre entidades linguísticas dotadas de força ilocucionária e, portanto, de autonomia pragmática e unidades linguísticas concluídas por uma quebra prosódica perceptível é um dado de fato. O WinPitch permitirá aos usuários verificar por conta própria a adequação da segmentação do fluxo da fala em enunciados proposta no C-ORAL-BRASIL, mas o recurso é acompanhado por uma série de medidas de validação e dos resultados em termos de acordo entre os anotadores, que deixam poucas dúvidas a respeito da sua confiabilidade. O C-ORAL-BRASIL dedicou uma parte importante das atividades de constituição do *corpus* à validação e ao refinamento dos resultados obtidos quanto à confiabilidade da partição do dado acústico em unidades discretas com base prosódica (veja-se o Capítulo 4 deste volume).

A proeminência das quebras prosódicas em geral foi verificada em muitas ocasiões. Quanto às pesquisas relativas à constituição dos grandes *corpora* de fala, o DUTCH *corpus*, no começo da década passada, apresentou dados relevantes com relação à detecção das quebras por parte de não *experts*. Especificamente a proeminência das quebras terminais para a percepção e, em medida consistente, também a proeminência das quebras não terminais, foram verificadas, sempre com não *experts*, nas coleções C-ORAL-ROM. Também foram replicadas como um dado relevante interlinguisticamente para além do quadro românico (vejam-se os trabalhos de Izre’el sobre o hebraico moderno (CoSIH) e a recente constituição do *corpus* das línguas afro-asiáticas (CorpAfroAs)). O trabalho desenvolvido para o PB não nasce, portanto, do nada, mas a estratégia de anotação e revisão adotada nesse caso é original e permitiu alcançar resultados qualitativamente muito superiores ao modelo do qual descende.

A percepção das quebras, de fato, não corresponde à percepção de um dado positivo, como, por exemplo, seria uma real interrupção do fluxo da



fala (pausa), mas a uma constelação de qualidades. *Reset* prosódico, queda de intensidade, proeminências prosódicas adjacentes e alongamento silábico são propriedades prosódicas, mas o cumprimento do ato de fala e a presença de fronteiras sintáticas são também correlatos que acompanham a percepção de uma quebra. Pode, portanto, ser difícil para um transcritor distinguir a informação prosódica relevante de outras informações, e isso pode levar tanto a fenômenos de sobre-extensão das marcas prosódicas na transcrição quanto à sua subextensão e, assim, resultar em uma diminuição da coerência das anotações. Raso e Mittmann (2009 e no Capítulo 4 deste livro) mostraram que essa coerência aumenta muito se os grupos de anotadores se comportam de maneira experiente e adequam seu horizonte perceptual, alcançando entre os anotadores um acordo nas fases anteriores à elaboração do *corpus*. Esta praxe deverá ser levada em conta na compilação de recursos futuros, e, em particular, a formação dos transcritores deverá ser considerada como um requisito essencial para a formação de *corpora* orais, se se quiser que os *corpora* se adêquem ao padrão de qualidade do C-ORAL-BRASIL.

A transcrição do C-ORAL-BRASIL contém sinais diacríticos de ruptura prosódica terminal e não terminal. O sistema é o mesmo aplicado a todos os *corpora* europeus do C-ORAL-ROM. Este requisito é central, a nosso ver, para o estudo da estrutura primária que caracteriza a organização da fala antes da sintaxe, ou seja, a articulação da informação. A escansão prosódica, que, repetimos, é um dado perceptual imediato e proeminente, é considerada na Teoria da Língua em Ato como a marca necessária da organização informacional com base pragmática dos enunciados e de entidades mais extensas da textualidade falada (padrões ilocucionários e estrofes). Sinteticamente, a arquitetura geral pode ser assim esquematizada:

- ato perlocucionário: ativação afetiva (pulsão) do ato da linguagem como um todo;
- ato ilocucionário: através do impulso afetivo, a ativação de um esquema acional segundo tipologias ilocucionárias convencionalizadas e a organização de um padrão informacional composto por várias unidades informacionais com diferentes funcionalidades pragmáticas, cujo núcleo deve ser constituído por uma unidade necessária e suficiente dedicada ao cumprimento da força ilocucionária do enunciado;



- interface prosódica obrigatória entre a sinalização do cumprimento do ato ilocucionário, as diversas unidades informacionais que podem compor o ato, e a sua execução locutiva;
- ato locutivo: através da sinalização prosódica, ativam-se entidades semântico/sintáticas correspondentes às unidades de informação componentes do padrão informacional e cumpre-se o ato ilocucionário, ou as outras unidades textuais maiores.

No Capítulo 3, Raso ilustra em detalhe as premissas teóricas e se aprofunda na descrição das diferentes funções informacionais, com uma rica exemplificação retirada do *corpus*. As funções informacionais podem ser classificadas segundo duas tipologias primárias: aquelas de tipo textual (Comentário, Tópico, Apêndice de Comentário e de Tópico, Parentético, Introdutor Locutivo) e aquelas de suporte dialógico (Incipitário, Fático, Conativo, Alocutivo, Expressivo, Conector Discursivo), que foram identificadas ao longo dos anos de trabalho experimental conduzido sobre o *corpus* italiano e sobre amostras espanholas e francesas, mas que, a esta altura, acharam confirmação na análise sistemática conduzida sobre o *corpus* de fala do PB. O C-ORAL-BRASIL dará, e já está dando, uma contribuição essencial ao estudo interlinguístico da articulação da informação na fala. Com base no C-ORAL-BRASIL e no C-ORAL-ROM italiano, foi possível constituir duas bases de dados comparáveis de fala informal de 5 mil enunciados por *corpus* etiquetados informacionalmente. No VII International LABLITA Workshop, realizado em Florença, em junho de 2011, foram apresentados os primeiros trabalhos de caráter geral sobre a estrutura informacional de línguas românicas como o italiano, o PB e o francês. Os colegas da equipe mineira começaram então a evidenciar peculiaridades da organização informacional, mas, conseqüentemente, também semântico-sintática, da variedade brasileira com relação às outras línguas românicas. O trabalho está só no começo, mas se anuncia como um dos pontos de maior novidade e relevância no estudo da fala.

Emanuela Cresti  
Massimo Moneglia  
(LABLITA, Università di Firenze)