

# *Desenvolvimento de um recurso léxico com papéis semânticos para o português*

*Leonardo Zilio (PPG-Letras/UFRGS)*

*Orientadora: Prof<sup>a</sup>. Dr<sup>a</sup>. Maria José Bocorny Finatto (PPG-Letras/UFRGS)*

*Coorientadora: Prof<sup>a</sup>. Dr<sup>a</sup>. Aline Villavicencio (PPGC/UFRGS)*

*Orientador no Exterior: Prof. Dr. Mathieu Mangeot-Nagata (LIG/UJF)*

*Coorientador no Exterior: Prof. Dr. Carlos Ramisch (LIF/AMU)*



# Informações Gerais

Uso de corpora especializado e não-especializado

Trabalho realizado no nível da oração

Anotação manual de informações semânticas

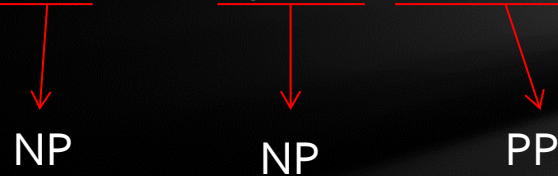
Essa anotação é feita com base em dois tipos de anotação já existentes

# Anotação de Subcategorização

João foi para casa.



João abriu a porta com a chave.

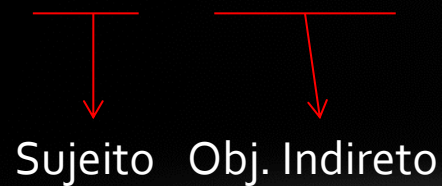


A porta abriu com a chave.



# Anotação Sintática

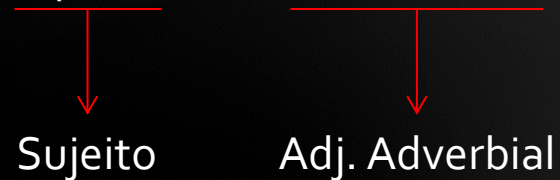
João foi para casa.



João abriu a porta com a chave.



A porta abriu com a chave.

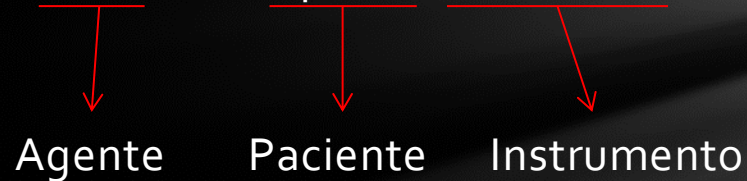


# Anotação de Papéis Semânticos

João foi para casa.



João abriu a porta com a chave.



A porta abriu com a chave.



# Objetivos

Gerar um recurso léxico com anotação de papéis semânticos

Contrastar as estruturas argumentais de verbos de linguagem especializada com os de linguagem não especializada

# Exemplo Especializado vs. Não Especializado

## Corpus de Cardiologia:

- "Atualmente *esse aparelho* pode ser *encontrado* nas unidades de *atendimento*, porém sua interpretação depende de especialistas, que muitas vezes não se encontram presentes em o momento de o exame."
- Estrutura: NP/tema + VP + PP/local

## Corpus DG:

- "*O pé direito do calçado* foi *encontrado* no buraco da loja de celulares."
- Estrutura: NP/tema + VP + PP/local

Os elementos linguísticos presentes são diferentes, mas a estrutura de papéis semânticos é a mesma.

*Corpora*



# *Corpus* de Cardiologia

Compilado por Zilio (2009)

490 artigos especializados extraídos de 3 periódicos  
brasileiros

> 1,5 milhão de palavras

Anotado com o *parser* PALAVRAS (Bick, 2000)

# Corpus Diário Gaúcho

Compilado no âmbito do projeto PorPopular

(<http://www6.ufrgs.br/textecc/porlexbras/porpopular/index.php>)

Textos retirados do jornal Diário Gaúcho (2008)

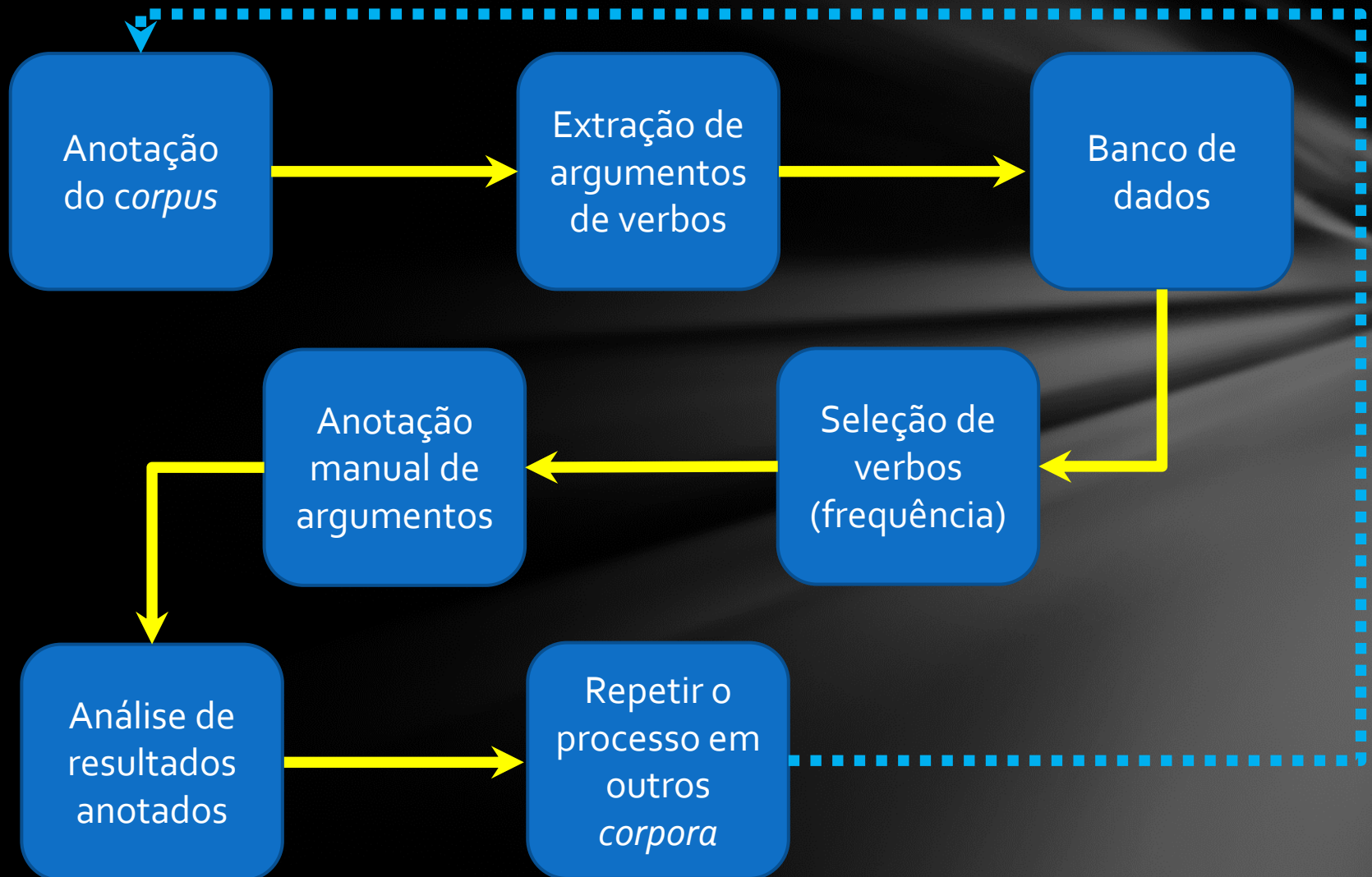
> 1 milhão de palavras

Anotado com o *parser* PALAVRAS (Bick, 2000)

# Metodologia

The background features a dark, almost black, field with several bright, white, diagonal light rays emanating from the right side, creating a sense of depth and focus. A solid, medium-green horizontal bar spans the bottom portion of the image.

# Etapas do Trabalho



# Extração automática de estruturas argumentais

Zanette (2010) - Extrator de *subcategorization frames* verbais a partir das árvores deitadas do PALAVRAS (BICK, 2000)

2011 = Modificação para árvore de dependências

Zanette et al. (2012) e Zilio et al. (2012) =  
Modificações nos parâmetros de reconhecimento de argumentos

# A anotação de papéis semânticos

## Exemplos do frame 'NP\_NP' do verbo 'apresentar'



Primeira « 1 2 » Última

### Exemplo 1

Dez pacientes apresentaram estimulação diafragmática , sendo o cabo-eletrodo reposicionado com sucesso em 9 pacientes .

Mostrar anotação

ARG_1	Dez pacientes	SUJEITO	Selezione	
ARG_2	estimulação diafragmática	OBJETO DIRETO	Selezione	
ARG_3	sendo o cabo-eletrodo reposicionado com sucesso em 9 pacientes	ADJUNTO ADVERBIAL	agente	

### Exemplo 2

Três pacientes apresentaram bloqueio atrioventricular total transitório após manipulação de o cabo-eletrodo em a cavidade ventricular direita , sendo necessária a estimulação ventricular provisória .

Mostrar anotação

ARG_1	Três pacientes	SUJEITO	dimensao geografica local	cione
ARG_2	bloqueio atrioventricular	OBJETO	local de origem	cione
ARG_3	total transitório após manipulação de o cabo-eletrodo em a cavidade ventricular direita sendo necessária a estimulação ventricular provisória	ADJUNTO ADVERBIAL	local de destino trajeto	cione

# Doutorado-Sanduiche

# Estudo-Piloto

50 verbos manualmente anotados nos  
dois *corpora*

Exportação dos dados para XML

Trabalho de Master I: Aluno Samy Sassi  
Conversão de MySQL para XML

Comparação das anotações nos dois  
*corpora*



# Dados da Anotação

3.400 Sentenças anotadas

1.790 de Cardiologia

1.610 do Diário Gaúcho

Cardiologia = 304 estruturas

Diário Gaúcho = 272 estruturas

# Estruturas mais frequentes nos dois corpora

<u>Diário Gaúcho</u>			<u>Cardiologia</u>		
Estrutura	Freq.	Freq. %	Estrutura	Freq.	Freq. %
SUJ<Agent>+OBJ.DIR<Theme>	171	10,62	SUJ<Theme>	181	10,11
SUJ<Theme>	114	7,08	SUJ<Theme>+ADJ.ADV[em] <Location>	121	6,76
SUJ<Agent>	92	5,71	SUJ<Instrument>+OBJ.DIR <Theme>	102	5,70
SUJ<Theme>+ADJ.ADV [em] <Location>	50	3,11	SUJ<Agent>+OBJ.DIR<Theme>	63	3,52
SUJ<Agent>+OBJ.DIR<Theme> +ADJ.ADV [em]<Location>	45	2,79	SUJ<Patient>	40	2,23

# Comparação Quantitativa

Coeficiente Tau-b de Kendall:

Correlação inexistente para estruturas sintático-semânticas ( $\tau = 0,031$ ).

SUBJ<Agent>+DIR.OBJ<Theme>

SUBJ<Experiencer>+DIR.OBJ <Theme>

Correlação existente para argumentos ( $\tau = 0,523$ ):

SUBJ<Agent>

SUBJ<Experiencer>

DIR.OBJ<Theme>

# Estudo-Piloto: Comparação Qualitativa

Maior uso de voz passiva e de intransitivos no *corpus* de Cardiologia

Maior presença de Instrumentos na posição de sujeito no *corpus* de Cardiologia

# Pra que serviu?

Aprimorar a lista de papéis semânticos para dar conta de mais tipos de argumentos

Ter uma ideia preliminar da diferença/semelhança entre as linguagens comum e especializada do ponto de vista dos papéis semânticos

# Atividades paralelas

Teste com vários anotadores

Transferência dos dados para a plataforma Jibiki

<http://jibiki.univ-savoie.fr/jibiki/Home.po>

# Teste com vários anotadores

10 anotadores linguistas com treinamento básico

25 sentenças extraídas do *corpus*

Cálculo de concordância entre anotadores (multi- $\pi$ )

Baixa concordância:

$$\pi = 0,253407$$

# Plataforma Jibiki

<http://jibiki.univ-savoie.fr/jibiki/Home.po>



# Plataforma Jibiki

Home Information Contacts Help

lang/语言/أَعَدُّ...

## Search result

USER: uebeltrager [User Profile](#) [Sign out](#)

1 entry(ies) retrieved.

LOOKUP:

Word: encontrar

Source: Portuguese

Target: All languages

Go

Advanced Lookup  
Dictionary List

ENTRIES:

Create  
Edit

REVIEW:

Contributions  
Alphabetical Lookup  
Contributors Board  
Export Volume

ADMINISTRATION:

Dictionaries  
Volumes  
Entries  
Stylesheets  
Users  
Groups  
Data Caches  
Edit news  
System Properties  
Server Statistics

FINISHED by [edit](#) [duplicate & edit](#) [delete](#) [view history](#) [view XML](#)

Verbo: encontrar Freq: 454

Estrutura: SUBJ\_V\_NP Voz: ATIVA Freq: 70

Estrutura: SUBJ\_V\_NP\_PP[em] Voz: ATIVA Freq: 34

Estrutura: SUBJ\_V\_PR Voz: PASSIVA Freq: 24

Estrutura: SUBJ\_V\_NP\_PR Voz: ATIVA Freq: 20

Estrutura: SUBJ\_V\_PP[em] Voz: PASSIVA Freq: 17

Estrutura: SUBJ\_V\_REFL Voz: ATIVA Freq: 16

Estrutura: SUBJ\_V\_PR\_PP[em] Voz: PASSIVA Freq: 14

Estrutura: SUBJ\_V\_REFL\_PP[em] Voz: ATIVA Freq: 12

Estrutura: SUBJ\_V\_NP\_PP[em]\_PP[em] Voz: ATIVA Freq: 12

Estrutura: SUBJ\_V Voz: PASSIVA Freq: 12

Estrutura: SUBJ\_V\_NP\_PP[a] Voz: ATIVA Freq: 8

Estrutura: SUBJ\_V\_REFL\_PP[com] Voz: ATIVA Freq: 7

Estrutura: SUBJ\_V\_PP[em] Voz: ATIVA Freq: 5

Estrutura: SUBJ\_V\_REFL\_PP[em]\_PP[em] Voz: ATIVA Freq: 5

# Plataforma Jibiki



Estrutura: SUBJ\_V\_NP    Voz: ATIVA    Freq: 70

- Encontrei um túmulo destruído , que não tinha dono , com os dois vasos .

## Argumentos

- *OCULTO*  
Sintaxe: SUJEITO    Papel Semântico: agente    Papel VerbNet:
  - *um túmulo destruído que não tinha dono com os dois vasos*  
Sintaxe: OBJETO DIRETO    Papel Semântico: tema    Papel VerbNet:
- Anderson beija Fernanda , quando encontra Klaus e Luiza .

## Argumentos

- *OCULTO*  
Sintaxe: SUJEITO    Papel Semântico:    Papel VerbNet:
  - *Klaus Luiza*  
Sintaxe: OBJETO DIRETO    Papel Semântico:    Papel VerbNet:
- Isso aqui é a melhor terapia que eu poderia ter encontrado .

## Argumentos

- *que*  
Sintaxe: OBJETO DIRETO    Papel Semântico:    Papel VerbNet:
  - *eu*  
Sintaxe: SUJEITO    Papel Semântico:    Papel VerbNet:
- Encontrei um bom acerto de o carro , gosto de essa pista e acho que seremos competitivos afirmou Felipe .

# Andamento

Modificação na lista de Papéis Semânticos

86 verbos anotados até agora

Modificação da interface da plataforma Jibiki

Qualificação do Doutorado

# Bibliografia

- BICK, Eckhardt. (2000) *The Parsing System PALAVRAS: Automatic Grammatical Analysis of Portuguese in a Constraint Grammar Framework*. Aarhus: Aarhus University Press. Disponível em: <http://beta.visl.sdu.dk/~eckhard/pdf/PLP20-amilo.ps.pdf>
- ZANETTE, Adriano. (2010) *Aquisição de Subcategorization Frames para Verbos da Língua Portuguesa*. Projeto de Diplomação. UFRGS. Orientadora: Aline Villavicencio.
- ZANETTE, A; SCARTON, C.; ZILIO, L. (2012) Automatic extraction of subcategorization frames from corpora: an approach to Portuguese. In: *Proceedings of PROPOR 2012 - Demonstration Session*. Coimbra, Portugal.
- ZILIO, Leonardo. (2009) *Colocações especializadas e Komposita: um estudo contrastivo alemão-português na área de Cardiologia*. Dissertação de Mestrado. Orientadora: Maria José Bocorny Finatto. Disponível em: <http://www.lume.ufrgs.br/bitstream/handle/10183/16877/000706196.pdf?sequence=1>
- ZILIO, Leonardo; ZANETTE, Adriano; SCARTON, Carolina (2012) Automatic extraction of subcategorization frames from corpora in portuguese. In: XI Encontro de Linguística de Corpus, São Carlos, SP, Brasil, 11-15 de setembro de 2012.

OBRIGADO!

Leonardo Zilio

[ziliotradutor@gmail.com](mailto:ziliotradutor@gmail.com)



# Papéis Semânticos/Temáticos/Theta ou Case Frames

Fillmore (1968) – Case Frame

Jackendoff (1987, apud Reinhart, 2002) – Relações Temáticas

Von Polenz (1988, apud Gelhausen, 2010) – Papéis Temáticos

Levin & Rappaport-Hovav (2005) – Papéis Semânticos

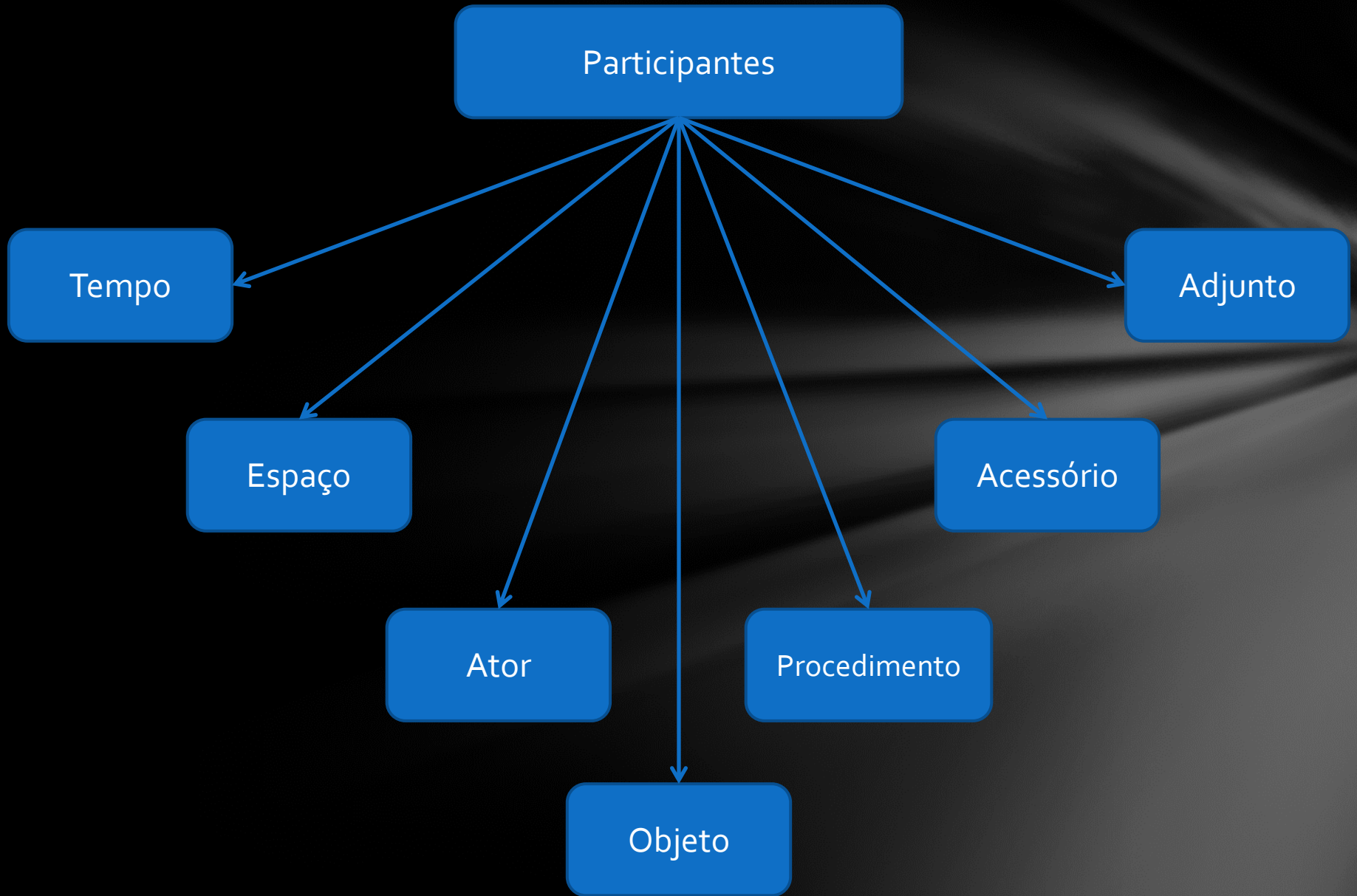
Gramática Gerativa – Theta Roles (Reinhart, 2002)

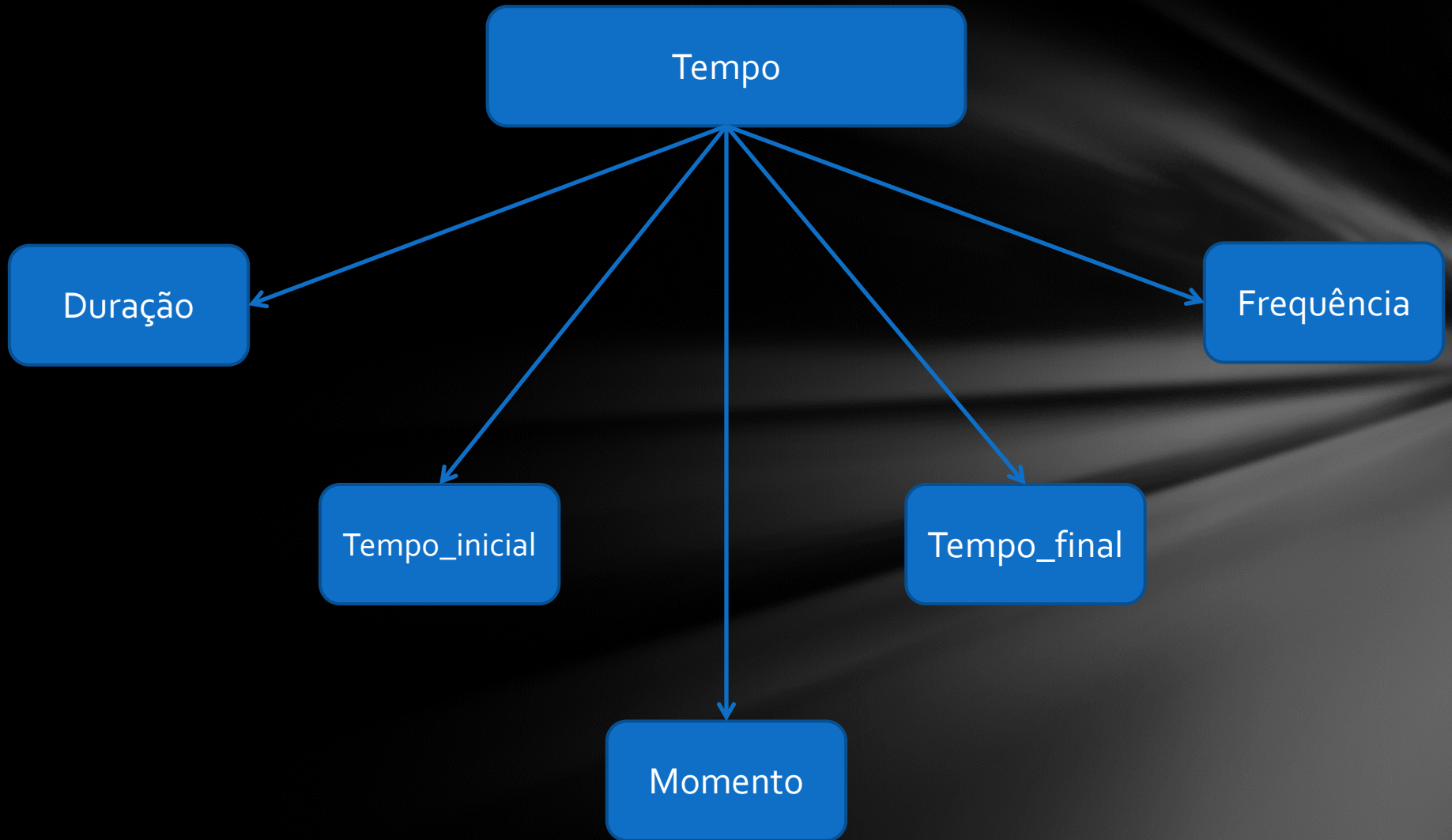
Dowty (1991) – Protopapéis Temáticos

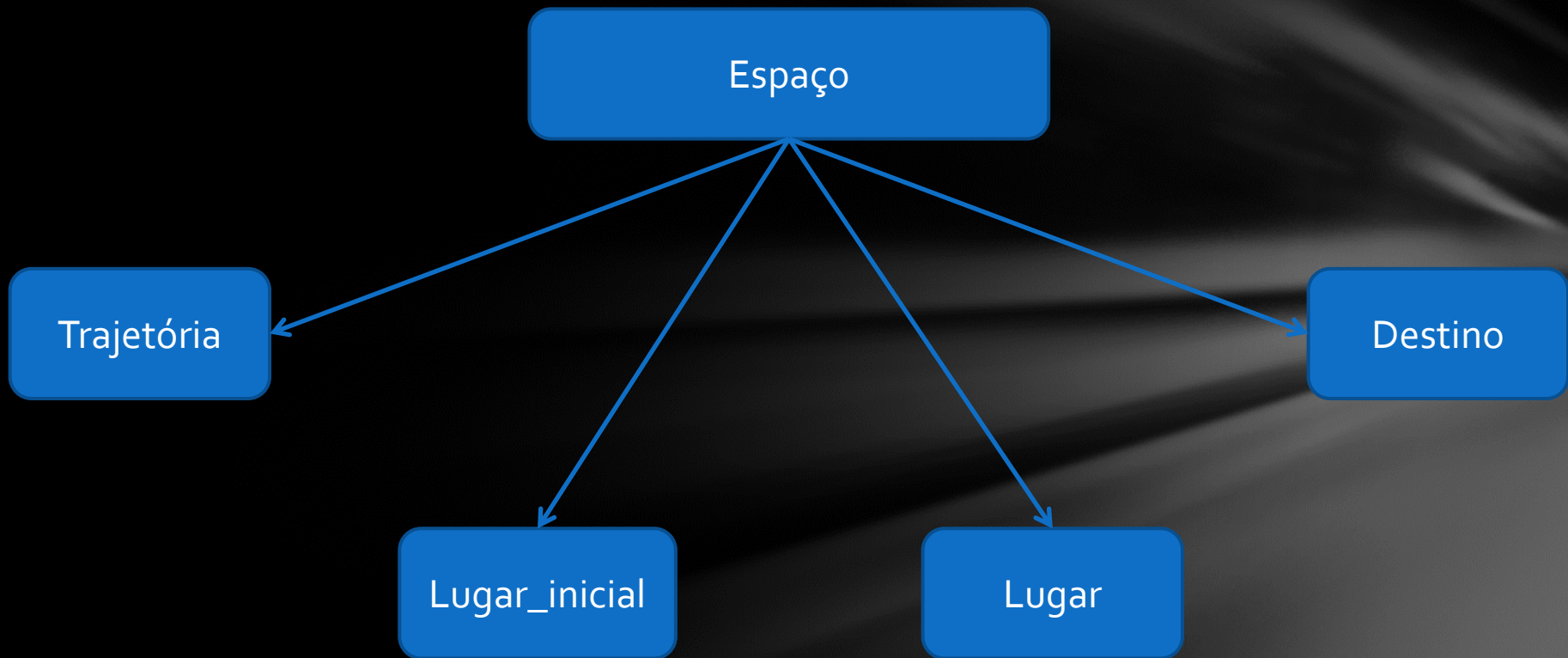
# Papéis Semânticos

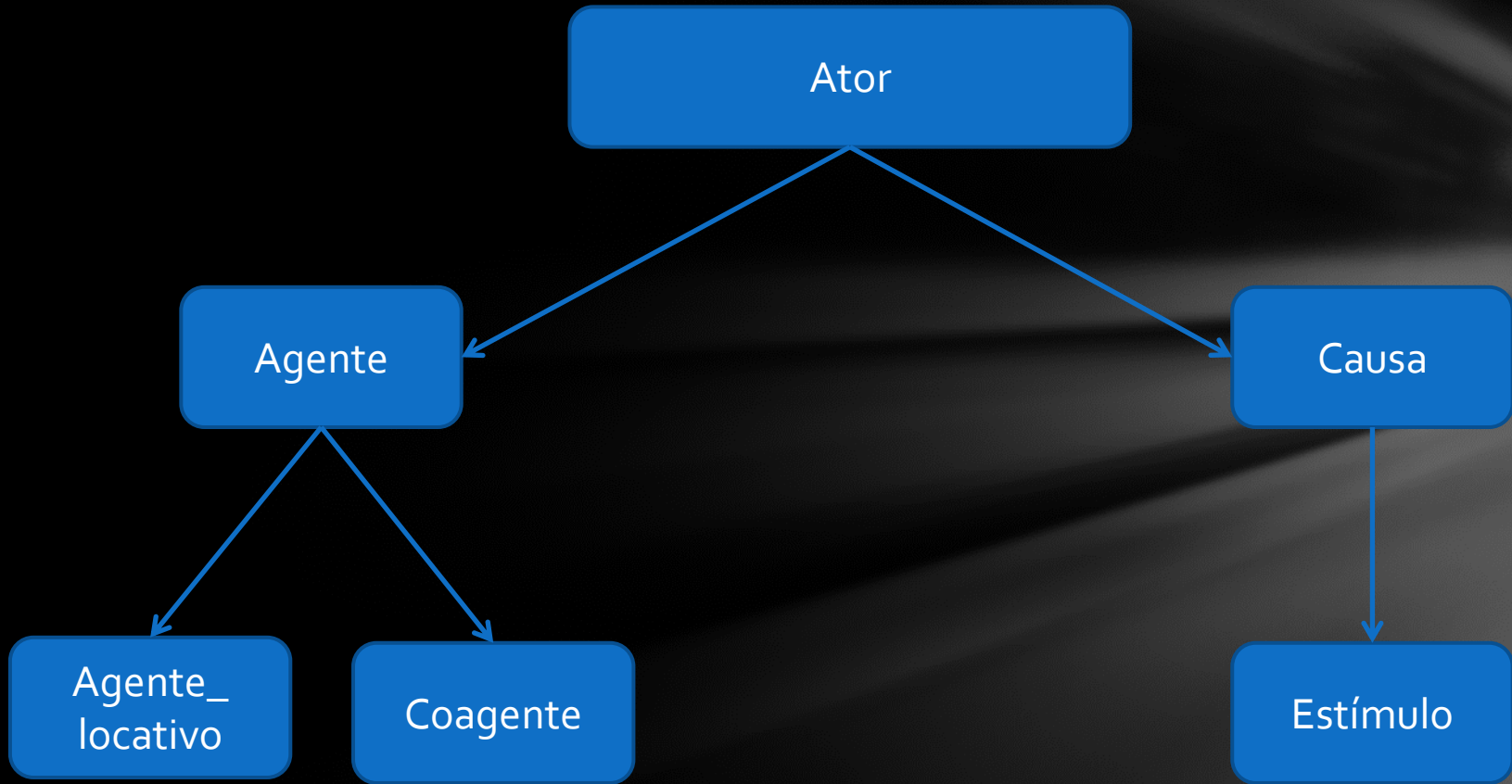
*Lista final estendida*











Objeto

Ter

Paciente

ente

Copaciente

Experien-  
ciador

