

UCC



Annotation of English

— Nathan Schneider ▫ COLING 2020 Tutorial: Part 2 —

<https://github.com/UniversalConceptualCognitiveAnnotation/tutorial>

Overview

Why UCCA?

- UCCA provides a **blueprint of conceptual compositionality** in a text—recognizing that sometimes
 - Semantic headedness ≠ syntactic headedness
 - Semantic predicate ≠ syntactic predicate (e.g. nouns can denote events)
 - Semantic “word”/minimal unit ≠ syntactic word (multiword expressions)
 - Semantic combinations may not be intuitively binary
 - A semantic dependent may be shared by multiple heads (syntax, inference)
 - Different languages use different grammatical trappings to convey information
- KEY DESIGN PRINCIPLES: **Foundational semantic graph structure, anchored in tokens, organized in terms of scenes, intuitive for annotators, multilingual, extensible with more layers**

Preliminaries

- Unit of annotation: **Passage**
- Base annotation layer: **Tokenization**
- This talk: the **Foundational Layer (FL)**
 - ▶ Main semantic graph structure in terms of **scenes**
 - ▶ FL depends on the tokenization + excluding punctuation
 - ▶ Other layers (e.g., **tense/aspect/modality, semantic roles, coreference**) can rest atop the FL

Preliminaries

- Unit of annotation: **Passage**
- Base annotation layer: **Tokenization**
- This talk: the **Foundational Layer (FL)**
 - Version 2 Guidelines dated November 24, 2018:
<https://github.com/UniversalConceptualCognitiveAnnotation/docs/blob/master/guidelines.pdf>

Foundational Layer

- Subsets of tokens are organized into **units**
 - Usually contiguous, but not always
 - A unit with 0 tokens is **implicit**
- Units nest within larger units
 - Minimal (terminal) units are disjoint in tokens, cover all non-punctuation tokens
 - * **Unanalyzable unit**: used for multiword expressions without internal head-modifier structure (e.g. person names, grammaticalized connectives)
 - Nonterminal units may contain other nonterminals and/or terminals
 - Annotation tool supports creating units top-down or bottom-up
- **Edges** relate units to parent units
 - Each edge contains a **category** label indicating the functional relationship of the child unit to the parent (e.g. Participant within a scene)
 - **Reentrancy**: If a unit has multiple parents (incoming edges), one is **primary**, others are **remote**
 - * Primary edges of a fully annotated passage always form a tree

Top-Level Units

- At the top level, the passage is segmented into units acting as Parallel Scenes (**H**) and **L**inkers

Sorkin conceived the political drama *The West Wing* in 1997 when he went unprepared to a lunch with producer John Wells and in a panic pitched to Wells a series centered on the senior staff of the White House, using leftover ideas from his script for *The American President*. He told Wells about his visits to the White House while doing research for *The American President*, and they found themselves discussing public service and the passion of the people who serve. Wells took the concept and pitched it to the NBC network, but was told to wait because the facts behind the Lewinsky scandal were breaking and there was concern that an audience would not be able to take a series about the White House seriously.

Top-Level Units

- At the top level, the passage is segmented into units acting as Parallel Scenes (**H**) and **L**inkers

Sorkin conceived the political drama The West Wing in 1997 when he went unprepared to a lunch with producer John Wells and in a panic pitched to Wells a series centered on the senior staff of the White House, using leftover ideas from his script for The American President. He told Wells about his visits to the White House while doing research for The American President, and they found themselves discussing public service and the passion of the people who serve. Wells took the concept and pitched it to the NBC network, but was told to wait because the facts behind the Lewinsky scandal were breaking and there was concern that an audience would not be able to take a series about the White House seriously.

Top-Level Units

- At the top level, the passage is segmented into units acting as Parallel Scenes (**H**) and **L**inkers

1 Sorkin conceived the political drama The West Wing in 1997 2 when he went
unprepared to a lunch with producer John Wells 3 and in a panic pitched to
Wells a series centered on the senior staff of the White House, using 4
leftover ideas from his script for The American President. He told Wells
about 5 his visits to the White House while doing research for The American
President, 6 and they found themselves discussing public service and the
passion of the people who serve. 7 Wells took the concept 8 and pitched it to
the NBC network, 9 but was told to wait 10 because the facts behind the
Lewinsky scandal were breaking and there was concern that an audience
would not be able to take a series about the White House seriously.

Top-Level Units

- At the top level, the passage is segmented into units acting as Parallel Scenes (**H**) and **L**inkers

Sorkin conceived the political drama The West Wing in 1997 when he went unprepared to a lunch with producer John Wells and in a panic pitched to Wells a series centered on the senior staff of the White House, using leftover ideas from his script for The American President. He told Wells [about [his visits to the White House] while [doing research for The American President]], and they found a common passion of the people with the NBC network, but was told to wait because the facts behind the Lewinsky scandal were breaking and there was concern that an audience would not be able to take a series about the White House seriously.

Parallel Scenes and Linkers may be embedded within a larger scene as well.

Top-Level Units

- At the top level, the passage is segmented into units acting as Parallel Scenes (**H**) and **L**inkers
 - [Josh started a fire]**H** **but****L** [unfortunately the chimney was blocked]**H**
 - **Either****L** [you come with me]**H** **or****L** [you stay at home]**H**
 - **After****L** [Abbey's party]**H** [we went to a bar]**H**



Scene Structure

- Main relation (scene-evoking unit): **S**tate or **P**rocess
- Participant (**A**) units
 - non-scene units (for most non-temporal NPs, PPs), as well as
 - scene units typically in a core syntactic position (subject, object/complement)
- Modifier units
 - **Adverbial** (**D**): manner/degree modifiers, modals, negation, spatial particles, ...
 - **T**ime: modifier (e.g. PP, adverb, adjective) expressing when or how often something happens *without constituting its own scene*
 - **G**round: extra-propositional element that relates a semantic unit to the speech event (speaker-oriented adverbial)

Main Relations: S vs. P

- **P**rocess: a dynamic event

- ▶ [Zoey presumably **graduates**_P from Georgetown tomorrow]_H
- ▶ [Zoey's **graduation**_P at Georgetown]_H
- ▶ cognitive activities like **seeing** & **thinking**: **P**

- **S**tate

- ▶ [Charlie really **loves**_S Zoey]_H
- ▶ [the block of cheese **weighed**_S 2 tons]_H



Participants (A)

- **P**rocess: a dynamic event

- ▶ [Zoey_A presumably **graduates**_P [from Georgetown]_A tomorrow]_H
- ▶ [[Zoey 's]_A **graduation**_P [at Georgetown]_A]_H
- ▶ cognitive activities like **seeing** & **thinking**: **P**

- **S**tate

- ▶ [[Charlie]_A really **loves**_S [Zoey]_A]_H
- ▶ [[the block of cheese]_A **weighed**_S [2 tons]_A]_H



Modifiers in Scenes

- **P**rocess: a dynamic event

- ▶ [Zoey_A presumably_G **graduates**_P [from Georgetown]_A tomorrow_T]_H
- ▶ [[Zoey 's]_A **graduation**_P [at Georgetown]_A]_H
- ▶ cognitive activities like **seeing** & **thinking**: **P**

- **S**tate

- ▶ [[Charlie]_A really_D **loves**_S [Zoey]_A]_H
- ▶ [[the block of cheese]_A **weighed**_S [2 tons]_A]_H



Participant vs. Adverbial

- Individuals, instruments, locations/destinations in an event are invariably **A**
 - [**Oliver**_A shattered_P **[the dictaphone]**_A **[with a hammer]**_A]_H
 - [**Leo**_A told_P **Bartlet**_A **[the news]**_A **[in his office]**_A]_H
- **D** applies only to units that do not introduce another participant or scene
 - [You_A **should**_D **not**_D behave_P **recklessly**_D]_H
 - [They_A treated_P him_A **[with disrespect]**_D]_H

Participant Scenes

- Scenes expressed with subjects, objects, and complement clauses can be **A**
 - **[the confirmation_P]_A** exhausted_P Toby_A]_H
 - [She_A announced_P **[that he had resigned_P]_A]_H**
 - [They_A broadcast_P **[her_A announcement_P [that he had resigned_P]_A]_A]_H**
- **scene unit** = any unit containing a **P** or **S** daughter.

Scenes: Practice

- Specify scene boundaries, Linkers, and each scene's main relation, Participants, and modifiers:
 - ▶ Jordan was annoyed when Leo angrily departed from the late meeting at the Capitol with Republicans



Scenes: Practice

- Specify scene boundaries, Linkers, and each scene's main relation, Participants, and modifiers:
 - ▶ Jordan was annoyed when Leo angrily departed from the late meeting at the Capitol with Republicans
 - ▶ [Jordan was **annoyed**_P]_H when_L [Leo angrily **departed**_P [from the late **meeting**_P at the Capitol with Republicans]]_H

Scenes: Practice

- Specify scene boundaries, Linkers, and each scene's main relation, Participants, and modifiers:
 - ▶ Jordan was annoyed when Leo angrily departed from the late meeting at the Capitol with Republicans
 - ▶ [Jordan was **annoyed**_P]_H when_L [Leo angrily **departed**_P [from the late **meeting**_P at the Capitol with Republicans]]_H
 - ▶ [Jordan_A was annoyed_P]_H when_L [Leo_A angrily departed_P [from the late meeting_P [at the Capitol]_A [with Republicans]_A]_A]_H

Scenes: Practice

- Specify scene boundaries, Linkers, and each scene's main relation, Participants, and modifiers:
 - ▶ Jordan was annoyed when Leo angrily departed from the late meeting at the Capitol with Republicans
 - ▶ [Jordan was **annoyed**_P]_H when_L [Leo angrily **departed**_P [from the late **meeting**_P at the Capitol with Republicans]]_H
 - ▶ [Jordan_A was annoyed_P]_H when_L [Leo_A angrily departed_P [from the late meeting_P [at the Capitol]_A [with Republicans]_A]_A]_H
 - ▶ [Jordan_A was annoyed_P]_H when_L [Leo_A angrily_D departed_P [from the late_T meeting_P [at the Capitol]_A [with Republicans]_A]_A]_H

Non-Scene Units

- If a non-scene unit has multiple children, the main one (semantic head): **C**enter
 - certain constructions warrant multiple Centers
 - **Q**uantity units
 - Connector (**N**) units
 - **E**laborator units
 - scene or non-scene
 - in general, **modifiers of non-scenes**: attributive adjective modifier, noun modifier in noun-noun compound, PP, apposition, relative clause, title, demonstrative determiner, degree modifier, ...
 - * (later: details on adjectives, appositions, relative clauses, PPs)
- [both presidential candidates_C with their wives]
- [all_Q 17_Q people_C]
- [Ed_C and_N Larry_C]
- [this_E chocolate_E cake_C]
- [Dr._E Bartlet_C]
- [Governor_C [of Maine]_E]
- [Lord_E [John Marbury]_C]
- [very_E angrily_C]

Functional Units

- Usually these are terminal units (no children)
- **R**elator units provide functional cues regarding a nested unit

- prepositions
- complementizers, relativizers: *that, which*
- subordinators that are not Linkers

[people_C [with_R hats_C]_E]
[plenty_Q [of_R hats_C]_C]
[He_A left_P [on_R Monday_C]_T]_H
[I_A saw_P [that_R he_A left_P]_A]_H

- **F**unction units

- articles
- non-modal auxiliaries
- copula with predicate adjective or relational noun
- expletive *it*
- polite forms
- infinitive *to* when not a purposive Linker

[the_F car_C]
[It_F will_F be_F raining_P]_H
[Could_F you_A please_F leave_P?]_H
[I waited [[for_R him_C]_A to_F leave_P]_A]_H

Lexical Units

- **Unanalyzable units (UNA):** multiple tokens forming a named entity or multiword expression where internal semantic structure is unclear. These multiword lexical units serve as leaves in the UCCA graph:
 - Personal names: *John Spencer*
 - Titles of works of art/literature/law: *The West Wing*
 - Dates: *May 17, 1832*
 - Foreign phrases: *Los Angeles, post hoc*
 - Idiomatic multiword expressions with opaque meanings: *hot dog, give up, in order to, as well as, according to, due to*
- **Generally analyzable:** proper names of places, organizations, and events, along with many specialized terms. Thus each token = 1 lexical unit.
 - University_C [of_R California_C]_E
 - 1600 Pennsylvania Ave., Washington, DC, USA – ?
 - time_E signature_C (in music)

Categories: Summary

<i>Parent unit:</i>	Scene	Non-scene only	Non-scene or Top-level
Required	P rocess or S tate	C enter	Parallel Scene (H)
Optional	Participant (A), Adverbial (D), T ime, G round	E laborator, Q uantity, Connector (N)	L inker
	Function, R elator		

Secondary categories:

UNanalyzable may be combined with any category in the table;
 Coordinated Main Relation (**CMR**) may occur with **P** or **S**

Basics: Practice

- Complete the parse:
 - [Jordan_A **was** annoyed_P]_H when_L [Leo_A angrily_D departed_P
[**from the** late_T meeting_P [**at the Capitol**]_A
[**with Republicans**]_A]_A]_H

Basics: Practice

- Complete the parse:
 - [Jordan_A **was**_F annoyed_P]_H when_L [Leo_A angrily_D departed_P
[**from**_R **the**_F late_T meeting_P [**at**_R **the**_F **Capitol**_C]_A
[**with**_R **Republicans**_C]_A]_A]_H

Basics: Practice

- Complete the parse:
 - [Jordan_A was_F annoyed_P]_H when_L [Leo_A angrily_D departed_P
[from_R [**the**_F]_{P-} late_T [**meeting**_C]_{-P} [at_R the_F Capitol_C]_A
[with_R Republicans_C]_A]_A]_H

Technically, determiners are attached to nouns within the main relation, creating a discontinuous unit [the_F meeting_C]_P

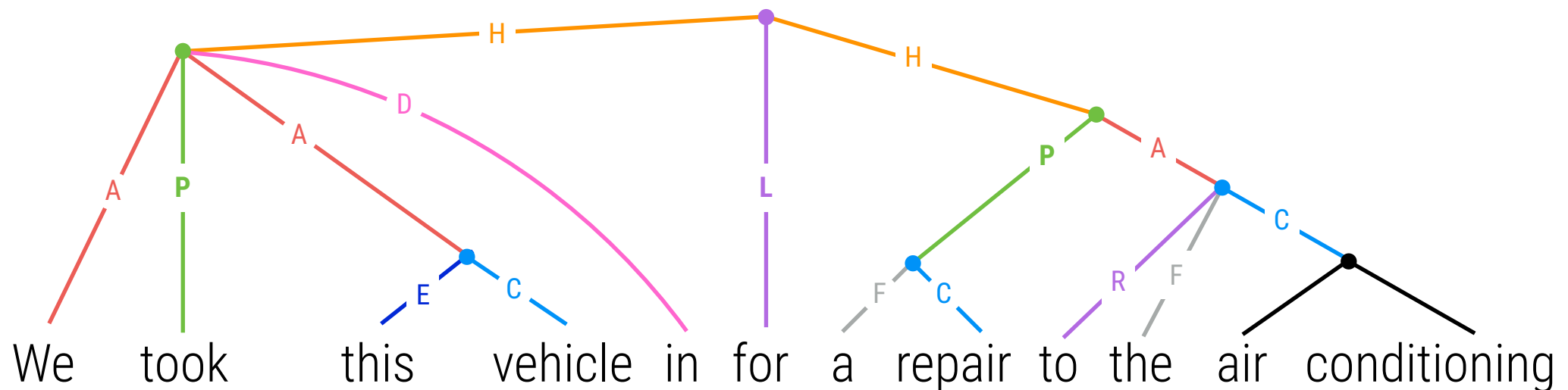
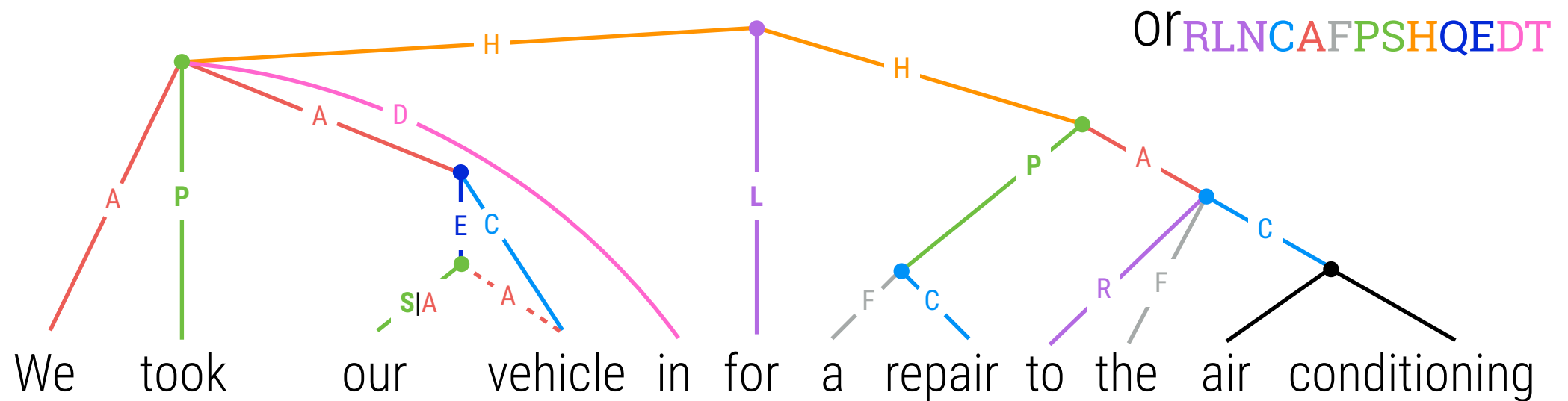
Simple Examples

1. Over the summer John read two books
2. Mary has been going to the gym every day for the last two years
3. John is speaking quietly and calmly to the children
4. Because we ate so early, we should bring a snack

Simple Examples

1. [[Over_R the_F summer_C]_T John_A read_P [two_Q books_C]_A]_H
2. [Mary_A has_F been_F going_P [to_R the_F gym_C]_A [every_E day_C]_T
[for_R the_F last_E two_Q years_C]_T]_H
3. [John_A is_F speaking_P [quietly_C and_N calmly_C]_D [to_R the_F
children_C]_A]_H
4. Because_L [we_A ate_P [so_E early_C]_T]_H, [we_A should_D bring_P
[a_F snack_C]_A]_H

Color scheme: **A D T G C E Q H P S N R L F**



We_A took_P this_E vehicle_C in_D for_L a_F repair_C to_R the_F [air conditioning]_C

We_A took_P this_E vehicle_C in_D for_L a_F repair_P [to_R the_F [air conditioning]_C]_A

We_A took_P this_E vehicle_C in_D for_L [a_F repair_C]_P [to_R the_F [air conditioning]_C]_A

We_A took_P this_E vehicle_C in_D for_L [[a_F repair_C]_P [to_R the_F [air conditioning]_C]_A]_H

[We_A took_P [this_E vehicle_C]_A in_D]_H for_L [[a_F repair_C]_P [to_R the_F [air conditioning]_C]_A]_H

English Constructions: A Tour

Adjectives / Remotes

- Predicative adjectives typically denote states:
 - [[the_F car_C]_A is_F red_S]_H
- Because many adjectives can be either attributive or predicative, we need a way to represent attributive adjectives as **both states and modifiers** simultaneously.
 - This is achieved with a **remote unit**, a reentrancy for the modified noun, denoted in parentheses: [I_A bought_P [the_F [red_S (car)_A]_E car_C]_A]_H
 - Thus the car token is shared between a non-scene unit in which it is the Center (its **primary edge**), and the scene unit where it is a Participant (a **remote edge**).
 - Every lexical unit has exactly 1 primary edge and may have 0 or more remote edges. We will see other uses of remote edges later.

Adjectives / Remotes

- However, not all uses of adjective are states: exceptions include
 - ▶ Pertainyms: a scientific_E paper_C
 - ▶ Modifiers of scene-evoking nouns: a beautiful_D wedding_P
 - * Only a non-scene unit can serve as Adverbial within a scene

CMR

- When predicate lexical units are **coordinated** it is tedious to annotate them as separate scenes with remote participants. A shorthand is to treat as non-scene coordination and mark the unit as **Coordinated Main Relation (CMR)**:
 - Adjectives: [He_A is_F **quiet_C and_N shy_C**_{S+CMR}]_H
 - Verbs: [Walden_A **wrote_C and_N recorded_C**_{P+CMR} [the score]_A]_H
- These can then be postprocessed to the full form:
 - [He_A is_F **quiet_S**_H **and_L** [(He)_A (is)_F **shy_S**]_H
 - [Walden_A **wrote_P** [the score]_A]_H **and_L** [(Walden)_A **recorded_P** (score)_A]_H

Degree Modifiers

Tough-Constructions

Secondary Predicates

- Depictives
- Resultatives

Role Nominals

- e.g. teacher

Raising & Control

- incl. Purpose clause control

Light/Secondary Verbs

Reflexives

Scenes within Scenes

- Complementation
- Subordination
- Relative Clauses
 - ▶ incl. preposition stranding
 - ▶ There is literallyG [no one [I do n't hate (no one)A [right now]T]E]

RCs: Technical Details

Questions

- wh; yes-no
- incl. free RCs?
 - Did you see what I did there ?



Copula Constructions

- ▶ $[CJ_A \text{ is } \mathbf{talls}]_H$
- ▶ $[[\text{The ambassador}]_A \text{ is } \mathbf{in}_S [\text{the Mural Room}]_A]_H$

Existentials

Possessives

Partitives / “of”

Appositions

- When two phrases are related by apposition,
 - ▶ if one of them is a name and the other is a description, the name is the Center and the description is the Elaborator
 - * Sheen portrays [[a fictional president]_E, [Josiah Bartlet]_C].
 - * Sheen portrays [[Josiah Bartlet]_C, [a fictional president]_E].
 - ▶ otherwise the syntactic head (usually the first item) is the Center and the other item is the Elaborator
 - * Sheen portrays [[the president]_C, [a Democrat]_E].

Implicit Units

- see: Infinitive clauses

Imperatives

Ellipsis

Comparative Constructions

Dates, Measurements, Ordinals

Vocatives, Interjections, Thanks

- also: direct quotations?

Fragments

More in Guidelines

- Question words
- Fused E-scenes
- Extraposition
- Directional particles
- Coordination
- Focus modifiers (“also”, “even”, “only”)
- Focus constructions (“It was ____ that ____”)
- **German** compound splitting
- ...

Summary

Formal Properties of Foundational Layer

Rooted DAG. Each edge has one or more **category** labels.

Primary Edges

- ▶ Form a tree (not necessarily projective in sentence order)
- ▶ Terminal unit: 0 or more non-punctuation tokens (2+ = unanalyzable unit); overt (non-implicit) units must be disjoint
- ▶ Units may be nested within other units, including unary nesting
- ▶ Units may be discontinuous
- ▶ Some simplifications are made prior to parser evaluation, which is span-based and forgiving w.r.t. attachment of F units (see later)

Remote Edges

- ▶ These are reentrancies within a passage (not necessarily same sentence)
- ▶ Can be grammatically required (e.g. control) or simply inferred
- ▶ Coreference between overt mentions (including pronouns) is NOT indicated in the foundational layer (but see section on extensions)

Categories

- cheatsheet

Acknowledgments