# QUD Structure as Discourse Structure: Segmentation, Labelling, and Genre Characteristics of Longer Texts

Oliver Deck        Tatjana Scheffler        Hannah Seemann

Ruhr University Bochum

## 1   Introduction

Questions under Discussion (QUD) are used to analyze discourse based on a stack of salient current questions that are seen as the basis of participants' dialogue moves (Roberts, 2012). In recent years, the framework has been used to investigate mainly short examples, lacking annotations of longer, natural texts, which may include false starts, disfluencies and other phenomena not seen in constructed discourse. The QUD-Anno Challenge encourages participants to annotate more substantial data to provide insights into the QUD framework's adaptability and applicability to different genres: a car review, a political interview, and a short fictional narrative.

In our annotation, we mainly focus on issues of text segmentation and QUD labelling that occur when working with longer, natural texts. We encountered challenges in the segmentation of verbal complements and conjoined phrases, disfluencies, and discourse markers. Such cases would profit from further clarification in the segmentation guidelines, but touch on the central question of what constitutes an "utterance" in discourse. Labelling and annotation challenges included the creative use of not-at-issue material (in the narrative text), the mismatch between implicit and explicit questions and answers (in the interview text), assertions starting with *but*, as well as general issues of formatting and numbering QUD tree structures.

The annotation was carried out by each of the three authors individually, following the annotation guidelines described in Riester et al. (2018). We focused on the narrative text and the interview; we have not yet annotated the car review. The individual annotation step was followed by a joint curation of only the segmentation. Based on this unified segmentation, we reannotated the QUD structure resulting in three separate annotations per text. Annotations based on the same text segmentation make further comparisons easier.

In the following section, we list some points that led to discussion when comparing our annotations.

## 2   Segmentation

The segmentation guidelines provided in Riester et al. (2018) are quite sparse, leading to a number of edge cases. While some of these cases could be resolved based on the annotation guidelines alone (like the example in Section 2.1), others required individual decisions based on discussion between the annotators. In general, more precise guidelines for how and when to split conjunctions, short sentences, etc., would be needed. This chapter lists issues that occurred during the segmentation of the data.

### 2.1   Conjoined Phrases

Conjoined phrases pose problems since it is not always clear whether or not the phrases should be considered (elliptical) clauses (i.e., segments). Following the guidelines of Riester et al. (2018), we decided to split the segments in (1) and decided in some cases to view conjoined noun phrases and even verb phrases are not separate assertions, as in (2).

(1)   Along the shelves are glass jars of coloured salts,
      and powders,
      and liquids.

(2)   But I think once we get it done, and once we can begin building a new partnership with our new friends, once we can start thinking about how we can do things differently, how we can interact with the rest of the world, how we can recover on our impetus, our mojos, as a global outward looking.

### 2.2   Disfluencies

Especially in the interview text, some utterances are not complete sentences, but incorporate false starts and repetitions. It might be argued, that these nevertheless convey meaning. For example, the turn in 3) can be seen as one utterance with a complex, left-extraposed subject NP. In our analysis of the discourse context, the first line in (3) turns out to be the answer to a preceding (implicit) question, as seen in (4). We thus annotated these utterances as though they were full assertions.

(3)   The extrication, after 45 years of our legal system, from the orbit of European law, which is you know, has become very, very pervasive.
      It's a very complicated thing to do.

(4)   A11: because obviously what the UK is going through is a big constitutional change.
      Q12: {What big constitutional change is the UK going through?}
      > A12: The extrication, after 45 years of our legal system, from the orbit of European law, which is you know, has become very, very pervasive.

> Q12.1: {How complex is the extrication from the orbit of European law?}
>> A12.1: It's a very complicated thing to do.

## 2.3 Discourse Markers

The interview text contained several sentences and clauses that served as discourse markers such as *Let's talk about that* and *So let's be completely clear* as seen in (5). We discussed whether these constitute their own assertions (especially when they have sentence form) or not. In this case, we decided to treat the discourse marker as part of the following segment. In general, the QUD model is meant to provide a structure of how discourse participants narrow down the space of possible worlds, which implies a language-external scope, i.e. not 'How do we want to communicate?' – 'completely clearly' but rather 'what are we communicating?'. We opted for integrating metacommunication such as the phrases in (5) into the surrounding segments whenever possible.

(5)     LK: Let's talk about that. So let's be completely clear, under the proposals that you were about to take to Brussels, there would be extra checks on the island of Ireland, how and where?

One should note that this may contradict the example in (Riester et al., 2018, pp. 13,20), who treats "We have all heard of conflicts" as a separate segment.

# 3 Labelling and Annotation

After creating the unified segmentation of both texts, we also encountered challenges in labelling/numbering and annotating the QUDs themselves; especially when assertions answered questions that were posed several utterances earlier or referred to content that was marked as non-at-issue. It is unclear whether the linear nature of the QUD stack adequately maps to more complicated, natural conversations as in the interview text or more creative, literary texts like the short story.

## 3.1 NAI in Narrative

In literary fiction, non-at-issue (NAI) content often drives the narrative. Artistic license allows meaning that is central to the story to be introduced in roundabout ways that must lead to inferences by the reader, as in the following example. In (6), *a play* is clearly the answer to an implicit question like *What does Susie see when she checks the calendar?*. However, in the previous sentence, *as she checks the calendar* is syntactically marked as not-at-issue, phrased in an embedded *as*-clause (cf. Potts, 2002). The at-issue clause is the second part of the sentence. It seems that NAI content can be used creatively in narrative to

advance the action and at the same time create a jarring effect for the reader. This is hard to reconcile with the QUD model, which depends on cooperation, and might thus not be a perfect fit for discourse that either has uncooperative participants (as in some parts of the interview text) or an author that flouts rules of conversational cooperation to create a literary effect.

(6)   As she checks Mrs Simpsons calendar, Susie rubs the place where the elastic cap from work scrunched all day.
A play.

## 3.2   Explicit vs. Implicit Questions

It is not always apparent whether a sentence or an utterance is a question or not. We discussed the different annotations in (7) and (8) and decided that the explicit turn (7–8) is not one explicit question, even though it ends with a question mark. We made this decision because the speaker does not stop to give space for an answer and the statement can be split up further into implicit questions. The analysis in (8) should therefore be preferred over the alternative shown in (7).

(7)   Q10: LK: But you're suggesting that people ought to come together, when transparently, you have been trying to create this idea of them and us,
> A10': you who want to get Brexit done, which you said every possible opportunity.
> A10": And the people on the other side, which you've just suggested, are only trying to hold you up and stop Brexit.
Q11: And that's transparent, you're trying to create a situation of them and us are you not?

(8)   Q7.1: {What is Boris Johnson saying about the people?}
> A7.1: LK: But you're suggesting that people ought to come together,
Q7.2: {What has Boris Johnson actually been doing to the public discourse?}
> A7.2: when transparently, you have been trying to create this idea of them and us,
> Q7.2.1: {What are the two sides of "them and us"?}
>> Q7.2.1.1: {Who is "us"?}
>>> A7.2.1.1: you who want to get Brexit done, which you said every possible opportunity.
>> Q7.2.1.2: {Who is "them"?}
>>>A7.2.1.2: And the people on the other side, which you've just suggested, are only trying to hold you up and stop Brexit.
Q7.3: And that's transparent, you're trying to create a situation of them and us are you not?

### 3.3 *But*

Assertions starting with *but* (very frequent in the narrative text) turned out to be difficult for annotation, because there is no question type that may be answered with a *but*-statement (cf. Scheffler, 2013, p. 88). For example, one cannot felicitously phrase (as an implicit question): "* What does this contrast with?", see (9).

(9)  Q3: {What is the tea like?}
     > A3: The tea is cold and bitter,
     Q4: {Does Susie care about the tea being cold and bitter?}
     > A4: but Susie doesn't mind

In cases like these, annotators must infer the question from the answer, since it is not possible to construct the question based on the previous sentences.

### 3.4 Side Stories/Backtracking

Complex narrative structure is difficult if not impossible to capture with the QUD stack and might call for a graph structure. Some implicit questions seem to repeat themselves, which cannot be annotated, even though it seems like the QUD model should be able to process such a "reopening" of questions. We used a solution shown in 1. This graph shows how all the questions of 'what does Susie do (next)?' are on the same graph level.

### 3.5 General Issues

Initial questions (especially in narrative) are often impossible to phrase without introducing new information. In (10), 'Susie' has to be introduced by the question. We were also uncertain on how to number questions and answers and did this more or less intuitively. More detailed guidelines would be helpful in this.

(10)  Q0: {What is the way things are?}
      Q1: {What does Susie do?}

## 4 Future Plans

In expanding this extended abstract, we first plan to evaluate our segmentation disagreements using standard inter-annotator agreement measures. We also want to evaluate our annotation agreement based on the curated (silver) segmentation, using QUD-specific tools[1], as well as tree-based F-scores and inter-annotator agreement.

---

[1]E.g., as provided in `https://github.com/QUD-comp/analysis-of-QUD-structures`.

```
├── [Q1.2 What does Susie do second?]
│   ├── As she checks Mrs Simpson's calendar,
│   ├── Susie rubs the place where the elastic cap from work scrunched all day.
│   ├── [Q7 What does Susie see?]
│   │   └── A play.
│   └── [Q8 What does Susie feel?]
│       ├── The skin feels puckered and soft, like white and wrinkly fingertips in the bath.
│       └── [Q9 What does that remind her of?]
│           ├── Her mother used to read beside the clawfoot bathtub her father imported from England.
│           └── [Q10 What would she do then?]
│               ├── Susie—who wasn't a Susie at all then—would tuck her chin over the edge of the tub and liste

                        ├── [Q11 What would happen to the water?]
                        │   └── The water would go cold.
                        └── [Q12 What would happen to Susie?]
                            └── Her skin would loose and crinkle.
├── [Q1.3 What does Susie do third?]
│   ├── Susie wipes the rim of the cup with her sleeve
│   ├── [Q1.3.1 Why is the tea fresh?] ## really a subquestion of Q4 or so
│   │   ├── Mrs Simpson only makes fresh tea for Mr Johnson next door.
│   │   └── [Q13 What is special about Mr Johnson's tea?]
│   │       └── One cup still contains a moss-smoke slick of whiskey.
│   ├── [Q14 What does the whiskey in the tea feel like?]
│   │   ├── The whiskey is too dilute to be warm
│   │   └── [Q14.1 Does Susie like the whiskey?]
│   │       └── but it's nice to know it's there.
│   └── then pours herself some tea.
├── [Q1.4 What does Susie do fourth after drinking the tea?]
│   ├── In the kitchen window, she sees the red elastic mark below her hairline.
│   └── [Q15 What reaction does she have?]
│       ├── She suppresses a sigh
│       ├── [Q15.1 What does she feel then?]
│       │   └── and the air sits heavy in her chest.
│       └── [Q16 What is the reason for the sigh?]
│           └── [Q17 How does Susie feel all day?]
│               ├── Susie has a knack for ignoring discomfit.
│               └── [Q18 How does that knack affect her life?]
│                   ├── It comes in handy at the factory.
│                   ├── [Q18.1 What is an example of that?]
│                   │   ├── She felt an itch that morning
│                   │   └── [Q19 What did she do about the itch?]
│                   │       └── but didn't scratch it until after her shift had ended.
│                   └── [Q18.2 What is  another example of enduring discomfit?]
│                       ├── She scrubs her face pink before each shift.
│                       ├── [Q20 Why is that necessary?]
│                       │   └── A single flake of dead skin can ruin a microchip.
│                       └── [Q21 How does that constitute a discomfit?]
│                           ├── But the factory is also vacuum-dry
│                           └── [Q22 What does she do against the dryness?]
│                               └── [Q22.1 What does she use against the dryness?]
│                                   ├── She can only afford sunflower oil.
│                                   └── [Q23 What happens when she moisturizes with sunflower oil?]
│                                       ├── For the whole day, she breathes the bitter staleness of cooking oil
```
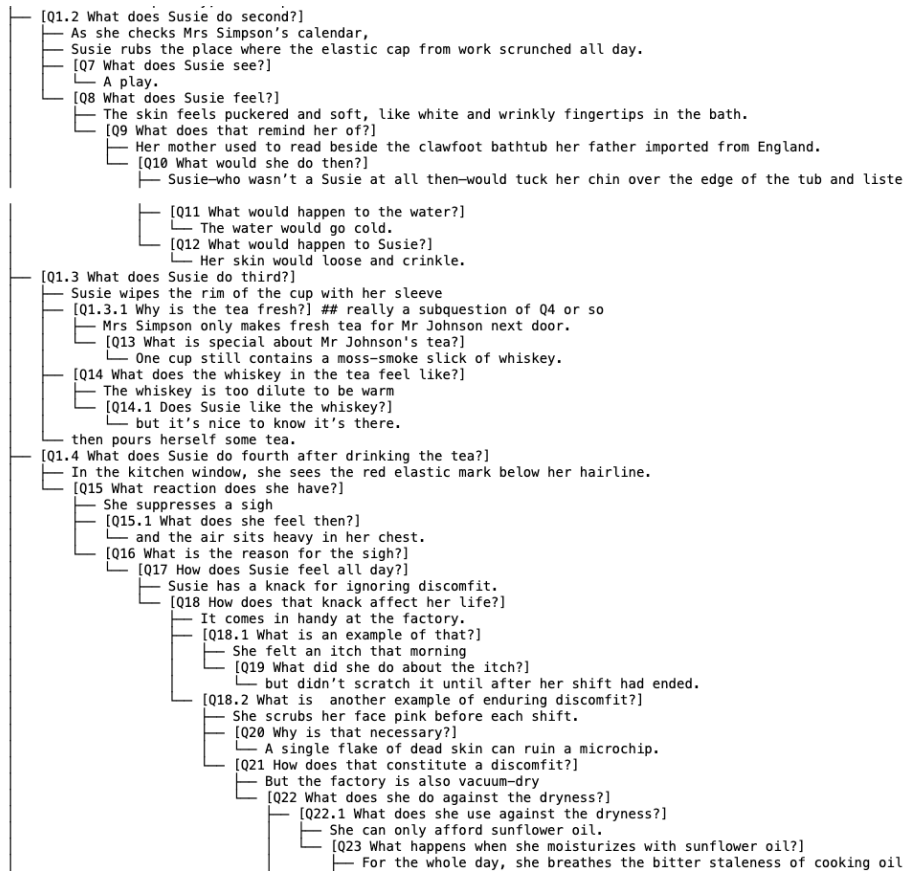
Figure 1: Visualization of part of the QUD tree structure of the narrative text.

We also look forward to providing (tentative) answers for the annotation issues mentioned above, including insights and input gained through discussion with the other workshop participants.

# References

Potts, Christopher. 2002. The syntax and semantics of as-parentheticals. *Natural Language & Linguistic Theory* 20. 623–689. doi:https://doi.org/10.1023/A:1015892718818.

Riester, Arndt, Lisa Brunetti & Kordula Kuthy. 2018. Annotation guidelines for questions under discussion and information structure. doi:10.1075/slcs.199.14rie.

Roberts, Craige. 2012. Information structure in discourse: Towards an integrated formal theory of pragmatics. *Semantics and Pragmatics* 5. 1–69. doi:10.3765/sp.5.6. http://semprag.org/article/view/sp.5.6.

Scheffler, Tatjana. 2013. *Two-dimensional semantics: Clausal adjuncts and complements*, vol. 549 Linguistische Arbeiten. Walter de Gruyter.