

0.1 What language? What grammar?

This proto-book is about aspects of the grammar of the contemporary Chinese language (现代汉语). Each word in this phrase can trigger controversy. Before we start substantial discussion, it is a wise idea to clarify what we are actually talking about.

0.1.1 Standard Mandarin and its variation

In the rest of the proto-book, we use the term *Chinese* 中文, (现代) 汉语, 华语 interchangeably with the more precise term *Standard Mandarin Chinese*.

0.1.2 The possibility to have a structure-based grammar

People familiar with Chinese often say its grammar depends more on the context, and some goes as far as claiming that a structure-based approach – or even a truth-value semantics-based approach – is infeasible when studying Chinese. And indeed, ?, arguably the most recognized grammar of the Chinese language, is a functionalist one. Our opinion is that though of course the context can influence strongly the grammar, this is not without limit. Pragmatic information may trigger *pro*-drop or forbid it, but it rarely triggers omission of the object. Sentences in dialogues may have more sentence final particles than written ones, but spoken sentences never have sentence *initial* particles. Though the context means a lot in Chinese, it is safe to assume we still have an underlying rigid structure beside purely semantic or pragmatic information.

The structure-based approach is by no means a rejection of functionalist studies. Rather, the former explores what features can be employed by the latter, so it can be expected the two approaches are complementary.

The next question is how to catch the structure. It is impossible to sit there and just “observe the world without bias”. People always do observation within a framework. Some may argue that typologists must be ready to invent completely new concepts when documenting a language (see, for example, ?), which is, of course, in principle true, but practically it is common for people to implicitly take some concepts for granted and carry out valuable works. R.M.W Dixon, a famous opponent of generative syntax, mocks “formalists” who fruitlessly try to find concepts exact corresponding to Indo-European ones in underdocumented languages in his *Basic Linguistic Theory* (BLT) (?) and advocates “describing a language in its own terms”, but he immediately goes on to discuss how to write a grammar *in terms of basic linguistic theory*, where we have predefined terms like *clause*, *sentence*, *argument*, a “deep structure” (This is indeed the term used by him!) which is made up by constituent hierarchy and so on.¹ Indeed, the strange fact that structuralist (and “arbitrary” and “purely empirical”) analyses of languages always fall into the same metalanguage – the one with headed (we will talk about the term in Section 0.2.3) phrase structures (IC analysis) and a set of shared concepts like predicate, arguments, etc. – is one motivation of the birth of generative syntax, which is formalized in Chomsky’s famous Syntactic Structures (?). We see the same fallacy in construction grammar, where people talk about stored routinized constructions – but routinization of *what*? It seems if we are to discuss purely structural aspects of a language, assuming a grammatical framework about possible structure building mechanisms is inevitable. This is actually not a bad thing. We will talk about the framework in following sections, and we will find Minimalism, tree-adjoining grammar, the implicit framework employed in many language documentation works, etc. can be reconciled.

0.1.3 “Not limited in Indo-European grammar studies”

Another frequently mentioned motto in the study of Chinese language is “Don’t be limited to Indo-European perspectives”. Again this is a correct statement but does not give much concrete

¹It is often justified that it is acceptable to do so because the predefined terms are just for inspiring people and not meant to be used in describing any language, and thus BLT is not a framework in the way generative approaches are. This justification is not valid, because in “hard sciences” like physics, it is quite common that a theory “breaks” the framework in a rigid sense but everyone agrees the theory is just *enriching* the framework, and there *is* a framework after all.

Table 1: Comparison between different formalisms

Feature	Minimalism	GB	TAG	Dependency grammar	BLT	CGEL
Surface-based segmentation	-	Descriptive works exist but not good	+	+	+	+
Fine grained atoms	+	-	-	-	-	-
Large “domains”	phase theory, cartography, etc.	+	+	-	+	+
Pre-compiled constructions	Nanosyntax-like lexicon		+	valency analysis	+	
Hierarchy details	+	+	+	-	-	+
Dependencies	Through DM-like subcategorization	Through notions like Spec-head relations	-	+	+	+

methodological suggestions.² What terms in Indo-European language studies should be avoided?

0.2 Existing descriptive framework

We list a few existing structure-based frameworks that are relevant for our discussions:

- Minimalism, which works on
- The traditional GB-style X-bar scheme
- **Tree-adjoining grammar (TAG)**
- Dependency grammars, which is frequently used in computational works.
- BLT, the framework used in most contemporary descriptive works.
- The grammatical framework used in The Cambridge Grammar of the English Language (CGEL) (??), which is generative-informed and yet remaining context-free and insists some analysis quite different from contemporary Minimalism (e.g. what is a head – we will discuss these apparent disagreements later).

Their differences can be roughly summarized as Table 1. In this section, we explain items in Table 1, the strength and weakness of each framework, and how these differences are mostly just notational differences and are more about methodology instead of worldview. The grammatical framework we use in this proto-book for descriptive works is a mixture of all these framework. We will also discuss why choosing such a framework and the relation between the framework and contemporary generative syntax.

²In self-identified “non-(or even anti-)generative” communities (the Language Hat, some Twitter circles, among others), this motto is also invoked to argue against formalist approaches. This accusation is very alarming and often contains many serious and insightful criticisms, but the claim itself may not factually hold, especially in recent years, since many generative linguistics are now highly interested in underdocumented languages, and many theoretical proposals (?) are based on these languages rather than so-called Indo-European perspectives. (Another related accusation is generative works do not view a language in a holistic way – how to solve the problem is also discussed in Section 0.2.1.) We should keep in mind that what works in English does not necessarily work in unfamiliar languages in question, but if a formal universal (for example, “the phonetic realization of pronouns is dependent to c-command relations”) seems truly reasonable in the new language, we should not hesitate to keep it.

0.2.1 Infeasibility of using derivational syntax as a descriptive tool

We first assess the “surface-based segmentation” and “find-grained atoms” row. Though Minimalism is the most prevalent framework in the generative enterprise, we should acknowledge that the framework is not a good choice for descriptive means (?), and other grammatical theories are required. The reasons can be summarized as follows:

- Minimalism tends to work with rather abstract features, which proves not suitable for language describing from sketch. When faced with real-world data, it is technically impossible for linguists to

0.2.2 Dependency relations

Dependency relations

0.2.3 Phrase structures

Divergent standards of constituency

Note that in the above discussion, a word – a bundle of several features spelt-out together – is a *span* and not a constituent in the generative sense, and yet we recognize it as a construction. This

So now we can see there is no conflict between the binary branching definition of *verb phrase* as Aux plus V plus NP_{object} and the BLT definition of verb phrase. A good choice is to use the term *verb complex* to denote verb phrase in the BLT sense (?).

There is yet a final issue about how flat the syntax tree should be.

The notion of *head*

The above discussion effectively gives us the good old X-bar theory, where a lexical word (not a functional word, not an invisible functional head) heads a phrase with multiple complements, specifiers and adjuncts (all in generative terms).

Types of dependents

It should be noted that the distinction between complements and modifiers are still opaque even with all these discussions about licensing and selection, and this opacity is theoretically rooted. Languages that have a relatively fixed word order often impose (though usually in a rather subtle way) a word order constraints on different types of clausal adjuncts and adjectives in NPs, which is well explained by the cartographic approach as we see above that the adjuncts and NP modifiers are actually introduced by a fixed hierarchy of functional heads. Since we consider the functional heads to be realized *onto* the CGEL head of the construction – the main noun, the main verb, etc. – it follows that adjectives about different properties fill different “slots of modifiers” of the head noun, and similarly adjuncts about different properties fill different “slots of modifiers” of the head verb, in exactly the same way arguments fill argument slots.

A clear distinction between complements and modifiers is therefore of limited purely formal interest and is better described as a language-specific concept, which may be useful for description (?).

Movements

Empty categories and fusion-head constructions

0.2.4 Categories and the notion of “word”

A **category** is defined as a type of constructions with similar distributions. We will first discuss basic syntactic constructions and identify positions that can be filled in them, and then search possible constructions in these positions. This is how categories can be recognized.

A construction that is small enough is said to be a **word** or more precisely, a **grammatical word**. We emphasize *grammatical* because it is quite common to use the term *word* 词 to denote a **prosody word**. Prosody is important in Chinese,

In traditional grammars concerning Latin, a common practice is to roughly define word classes (nouns, verbs, etc.) according to their meaning and then discuss where they can be used. In this proto-book we do not take this approach. Though we will review a lot of work based on the meaning-first approach, the way we distinguish word classes is mainly distributional. If two words can appear in similar positions, they are classified into one **word class** or **part of speech**. A word class is just a category about words.

In other words, we define concepts like *noun-like* and *verb-like* before listing criteria of what is a noun and what is a verb. Criteria for word classes are always language-specific, but we have more confidence that at least some *features* – like the nominal feature *n* or the verbal feature *v* – are cross-linguistic and may be attributed to the language faculty in the broad sense.

0.3 The mixed framework

After length discussions, it is time to go back and summarize the framework we adopt in this proto-book.

Having worked out a *formal* grammatical framework, we still need to go through a list of constructions that are

0.4 Data sources