**Presentation Abstract**

# Centaurus: A Just-in-time Parallel-parser Generator for Ad Hoc Data Processing

SHIGEYUKI SATO[1,a)]   HIROKA IHARA[2]   KENJIRO TAURA[1,b)]

Presented: July 26, 2019

It is important to handle data in text formats such as XML, JSON, and CSV because these data very often appear in the context of data exchange. Only parts of these data are typically used afterwards so that it is not worth ingesting the whole of them into databases. It is therefore desired to match and extract the concerned part in a lightweight ad hoc manner. Classically used for such a purpose are linewise regular expression tools such as grep, sed, and awk. These are, however, not powerful enough for text formats commonly used for data exchange because they cannot recognize nested structures in general. To support a lightweight ad hoc data processing, we present Centaurus, a just-in-time parallel-parser generator library. By generating native scannerless LL(*) parsers dynamically, our library enables us to process input data in parallel merely by calling Python functions with LL(*) grammars and Python actions. This presentation gives the design and implementation of Centaurus and reports its experimental performance on data filtering.