

# Hybrid Word-Subword Speech Recognition – a Powerful Tool to Search in Speech

*Jan ČERNOCKÝ, Igor SZOKE, Mirko HANNEMAN, Stefan KOMBRINK, Michal FAPŠO*

Brno University of Technology, Faculty of Information Technology, Božetěchova 2, 612 66 Brno, Czech Republic

{cernocky,szoke,ihannema,kombrink,ifapso}@fit.vutbr.cz

**Abstract.** Main-stream systems for searching information in speech are based on Large Vocabulary Continuous Speech Recognizer (LVCSR) with fixed vocabulary. The keywords or key-phrases are subsequently searched in its output. These systems have severe problems with Out of Vocabulary (OOV) words, that are common when one changes the domain (for example from standard to medical), speaker (normal versus highly educated), or even date (new words appearing in TV news). This talk will present our work in designing hybrid word-subword recognition systems, that have a combined recognition network. Under normal circumstances, they output standard word strings, while they are allowed to switch to subword description for unknown inputs. Such systems are good not only for detecting OOVs, but also subsequent steps leading to their exploitation. Under the EC-sponsored DIRAC project, we have investigated analysis of detected OOVs, conversion to standard word-form, and finding links to in-vocabulary words and other OOVs. The results will be demonstrated on real speech data from popular TED lectures.