



Bahnaric

Dơ gum lư bơ gloh 50% lai xuất ăn rim khoan xre wã pom
hnam, thiết 'bĩ Bình Minh lôm 2 xơ năm pơ tơm chong bi
gơh h'oh 500 triệu đơng lôm minh kơ sơ Đà Nẵng.

Loanword removal

Đor gum lư bơ gloh lai xuất ăn rim khoan xre wã pòm
hnam thiết 'bĩ lòm xơ năm pơ tòm chong bi gơh h'oh triểu
đông lòm minh kơ sơ

"Hỗ trợ", "lãi xuất", "cho", "các", "khoản vay",
"đề", "đầu tư", "nhà", "thiết bị", "trong"

"lu_bơ_gloh", "xơ_năm_pơ", "tôm_chong", "bì",
"goh h'oh", "riều_đông", "lơm", "minh", "kơ_sơ"

Unmapped phrases

Mapped phrases



Vietnamese

Our Hybrid NMT Architecture

Loanword Detection



Underthesea

Named entities: "Bình Minh", "Đà Nẵng"

Numerical values: "50", "2", "500"

Punctuation: "%", ",", ";", "."

Word Segmentation

$$\text{PMI}(x_1, x_2, \dots, x_n) = \log_2 \left(\frac{P(x_1, x_2, \dots, x_n)}{P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_n)} \right)$$

"dơ_gum", "lư_bơ_gloh", "lai_xuất", "ăn", "rim",
"khoan_xre", "vã", "pơm", "hnam", "thiết_bì",
"lơm", "xơ_năm_pơ", "tơm_chơng", "bì",
"gơh h'oh", "riều_đơng", "lơm", "minh", "kơ_sơ"

Lexical Mapping



Bahnaric-Vietnamese bilingual dictionary

BARTBahnar

"tối đa", "năm đầu", "nhưng", "không", "vượt
quá", "triệu đồng", "trên", "một", "cơ sở"

Post-Processing