# Generating a Parallel Corpus Stream for Odia: Mining Parallel Corpus from Odia Twitter

**Anjan Kumar Panda, MSC IT, KSOU Mysore**
Internet Application Specialist, Technology Manager
Life Member, OSA. The Odisha Society of the Americas
5050, Hacienda Drive, Apt 2232, Dublin, CA, 94568
panda.anjankumar@gmail.com
Contact: 1- 845-535-0961

**Dr Arun Kumar Malik, PhD, Assistant Professor of Political Science**
Gujarat National Law University, Gandhinagar
amalik@gnlu.ac.in
Contact No. 8128650850

Note: This research article is part of a sequence of papers to enable the researchers in the field of computational linguistics in Odia.

==================================================================
**Language in India** www.languageinindia.com **ISSN 1930-2940 20:8 August 2020**
Anjan Kumar Panda, MSC IT, KSOU Mysore and Dr Arun Kumar Malik, PhD
Generating a Parallel Corpus Stream for Odia: Mining Parallel Corpus from Odia Twitter      1

**Introduction**

A corpus is a fundamental need for natural language process applications.

A parallel corpus is a foundational need for languages like Odia (Oriya - The Unicode Standard, Version 13.0."https://unicode.org/charts/PDF/U0B00.pdf. Accessed 8 Aug. 2020) which would enable explorations in natural language processing advancements into machine translation, (Machine translation - Wikipedia. "https://en.wikipedia.org/wiki/Machine translation. Accessed 8 Aug. 2020) computational language modelling, (Language model - Wikipedia."

===============================================================
**Language in India** www.languageinindia.com **ISSN 1930-2940 20:8 August 2020**
Anjan Kumar Panda, MSC IT, KSOU Mysore and Dr Arun Kumar Malik, PhD
Generating a Parallel Corpus Stream for Odia: Mining Parallel Corpus from Odia Twitter   2