![turnitin](turnitin logo)

# Digital Receipt

This receipt acknowledges that Turnitin received your paper. Below you will find the receipt information regarding your submission.

The first page of your submissions is displayed below.

| | |
|---|---|
| Submission author: | Ye Kyaw Thu |
| Assignment title: | NLP Paper |
| Submission title: | myTTS-ver3 |
| File name: | nlp-proj.pdf |
| File size: | 200.7K |
| Page count: | 4 |
| Word count: | 2,080 |
| Character count: | 11,160 |
| Submission date: | 28-Jan-2021 06:03PM (UTC+0630) |
| Submission ID: | 1496175436 |