# 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: | SCSE STYBORSKI JEREMY ANDREW |
| Assignment title: | EE6222 Assignment: Action Recognition in the Dark |
| Submission title: | Styborski_DarkHAR_Report |
| File name: | Styborski_DarkHAR_Report.pdf |
| File size: | 292.25K |
| Page count: | 5 |
| Word count: | 3,836 |
| Character count: | 20,855 |
| Submission date: | 03-Nov-2022 11:57AM (UTC+0800) |
| Submission ID: | 1943101923 |



Deep Learning 3D CNN Applied to Low-Light Human Action Recognition Videos