# Deep learning for visual understanding

Guo, Y.; Guo Y.

**Citation**

Guo, Y. (2017, October 5). *Deep learning for visual understanding.* Retrieved from https://hdl.handle.net/1887/52990

| Version: | Not Applicable (or Unknown) |
|---|---|
| License: | [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#) |
| Downloaded from: | [https://hdl.handle.net/1887/52990](https://hdl.handle.net/1887/52990) |

**Note:** To cite this publication please use the final published version (if applicable).

Cover Page





The handle http://hdl.handle.net/1887/52990 holds various files of this Leiden University dissertation.

**Author**: Guo, Y.
**Title**: Deep learning for visual understanding
**Issue Date**: 2017-10-05

# Deep Learning for Visual Understanding

**Proefschrift**

ter verkrijging van
de graad van Doctor aan de Universiteit Leiden
op gezag van Rector Magnificus prof.mr. C.J.J.M. Stolker,
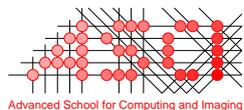volgens besluit van het College voor Promoties
te verdedigen op donderdag 5 oktober 2017
klokke 10.00 uur

door

## Yanming Guo

geboren te Hebei, China
in 1989

**Promotiecommissie**

Promotor:        Prof. Dr. J.N. Kok
Co-promotor:     Dr. M.S. Lew
Overige leden:   Prof. Dr. A. Plaat
                 Prof. Dr. W. Kraaij
                 Prof. Dr. T.H.W. Bäck
                 Prof. Dr. H. Trautmann (University of Münster, Germany)
                 Prof. Dr. S. Rüger        (Open University, United Kingdom)



This work was carried out in the ASCI graduate school.
ASCI dissertation series number: 378

# Contents