# 2025 IEEE Hot Chips 37 Symposium (HCS 2025)

**Stanford, California, USA**
**24-26 August 2025**

**Additional Copies of This Publication Are Available From:**

# TABLE OF CONTENTS

**Author Index**