# 37th Pacific Asia Conference on Language, Information and Computation (PACLIC 37)

Hong Kong, China
2-4 December 2023

# TABLE OF CONTENTS

**Author Index**