# 9th International Conference on Data Science, Technology and Applications (DATA 2020)

Online
7 – 9 July 2020

**Editors:**

**Slimane Hammoudi**
**Christoph Quix**
**Jorge Bernardino**

# CONTENTS

## SHORT PAPERS