Balázs Vőneki and Sébastien Valat

# Evaluation of 100 Gb/s LAN networks for the LHCb DAQ upgrade

Balázs Vőneki and Sébastien Valat

**Poster 229**

# Evaluation of 100 Gb/s LAN networks for the LHCb DAQ upgrade

Balázs Vőneki and Sébastien Valat

## Poster 229



Introduction

# Evaluation of 100 Gb/s LAN networks for the LHCb DAQ upgrade

Balázs Vőneki and Sébastien Valat

**Poster 229**

Introduction →

← Architecture

# Evaluation of 100 Gb/s LAN networks for the LHCb DAQ upgrade

Balázs Vőneki and Sébastien Valat

**Poster 229**



Introduction

Architecture

Various benchmark results