**PSL★**
UNIVERSITÉ PARIS

Préparée à MINES ParisTech

**THÈSE DE DOCTORAT**
DE L'UNIVERSITÉ PSL

# Network-guided genome-wide association studies

## Études d'association génome entier guidées par des réseaux

Soutenue par
**Héctor Climente González**
Le 4 Février 2020

Fin de confidentialité
Le 4 Février 2021

École doctorale n°621
**Ingénierie des Systèmes, Matériaux, Mécanique, Énergétique**

Spécialité
**Bio-informatique**

Composition du jury :

| | | |
|---|---|---|
| Nadine ANDRIEU<br>Mme., Institut Curie | | *Présidente* |
| Kristel VAN STEEN<br>Mme., Université de Liège | | *Rapporteuse* |
| Antonio RAUSELL<br>M., Imagine Institute | | *Rapporteur* |
| Laura FURLONG<br>Mme., Pompeu Fabra University | | *Examinatrice* |
| Véronique STOVEN<br>Mme., MINES ParisTech | | *Directrice de thèse* |
| Chloé-Agathe AZENCOTT<br>Mme., MINES ParisTech | | *Co-encadrante* |

# Contents