

# SPiP: Splicing Prediction Pipeline, a machine learning tool for massive detection of exonic and intronic variant effect on mRNA splicing.

Raphaël Leman<sup>1</sup>, Béatrice Parfait<sup>2</sup>, Dominique Vidaud<sup>2</sup>, Emmanuelle Girodon<sup>2</sup>, Laurence Pacot<sup>2</sup>, Gérald LE GAC<sup>3</sup>, Chandran Ka<sup>3</sup>, Claude Ferec<sup>3</sup>, Yann Fichou<sup>3</sup>, Céline Quesnelle<sup>1</sup>, Camille Aucouturier<sup>1</sup>, Etienne Muller<sup>1</sup>, Dominique Vaur<sup>1</sup>, Laurent Castera<sup>1</sup>, Flavie Boulouard<sup>1</sup>, Agathe Ricou<sup>1</sup>, Hélène Tubeuf<sup>4</sup>, Omar Soukarieh<sup>4</sup>, Pascaline Gaildrat<sup>4</sup>, Florence Riant<sup>5</sup>, Marine Guillaud-Bataille<sup>6</sup>, Sandrine Caputo<sup>7</sup>, Virginie Moncoubier<sup>7</sup>, Nadia Boutry-Kryza<sup>8</sup>, Françoise Bonnet-Dorion<sup>9</sup>, Ines Schultz<sup>10</sup>, Maria Rossing<sup>11</sup>, Olivier Quenez<sup>4</sup>, Louis Goldenberg<sup>4</sup>, Valentin Harter<sup>1</sup>, Michael Parsons<sup>12</sup>, Amanda Spurdle<sup>12</sup>, Thierry Frébourg<sup>4</sup>, Alexandra Martins<sup>4</sup>, Claude Houdayer<sup>4</sup>, and Sophie Krieger<sup>1</sup>

<sup>1</sup>Centre Francois Baclesse Centre de Lutte Contre le Cancer

<sup>2</sup>Hopital Cochin Service de Radiologie

<sup>3</sup>Universite de Bretagne Occidentale

<sup>4</sup>Universite de Rouen

<sup>5</sup>Groupe Hospitalier Saint-Louis Lariboisiere et Fernand-Widal

<sup>6</sup>Gustave Roussy

<sup>7</sup>Institut Curie Departement d'Oncologie Medicale

<sup>8</sup>Hospices Civils de Lyon

<sup>9</sup>Institut Bergonie

<sup>10</sup>Centre Paul Strauss

<sup>11</sup>Rigshospitalet

<sup>12</sup>QIMR Berghofer Department of Genetics and Computational Biology

February 21, 2022

## Abstract

Modeling splicing is essential for tackling the challenge of variant interpretation as each nucleotide variation can be pathogenic by affecting pre-mRNA splicing *via* disruption/creation of splicing motifs such as 5'/3' splice sites, branch sites or splicing regulatory elements. Unfortunately, most *in silico* tools focus on a specific type of splicing motif, which is why we developed the Splicing Prediction Pipeline (SPiP) to perform, in one single bioinformatic analysis based on machine learning approach, comprehensive assessment of variant effect on different splicing motifs. We gathered a curated set of 4,616 variants scattered all along the sequence of 227 genes, with their corresponding splicing studies. Bayesian analysis provided us the number of control variants, *i.e.* variants without impact on splicing, to mimic the deluge of variants from high throughput sequencing data. Results show that SPiP can deal with the diversity of splicing alterations, with 83.13% sensitivity and 99% specificity to detect spliceogenic variants. Overall performance as measured by area under the receiving operator curve was 0.986, significantly better than 0.965 spliceAI for the same dataset. SPiP lends itself to a unique suite for comprehensive prediction of spliceogenicity in the genomic medicine era. SPiP is available at: [<https://sourceforge.net/projects/splicing-prediction-pipeline/>] (<https://sourceforge.net/projects/splicing-prediction-pipeline/>)

## Hosted file

article\_SPiP\_manuscript\_vf.docx available at <https://authorea.com/users/461595/articles/557263-spip-splicing-prediction-pipeline-a-machine-learning-tool-for-massive-detection-of-exonic-and-intronic-variant-effect-on-mrna-splicing>



