# Convergence Analysis of a Momentum Algorithm with Adaptive Step Size for Nonconvex Optimization

**Anas Barakat**

**Joint work with Pascal Bianchi**

**LTCI, Télécom Paris, Institut polytechnique de Paris**

# A Momentum Algorithm with Adaptive Step Size

- ▶ ADAM famous **BUT** convergence issues (Reddi et al., 2018).

- ▶ Several variants : Yogi, AdaBound, AdaShift, Nadam, QHAdam, RAdam ...

- ▶ **Goal :** convergence rates for adaptive algorithms (ADAM in particular) for **nonconvex** optimization.

## Algorithm

$$
\begin{cases}
x_{n+1} = x_n - a_{n+1} p_{n+1} \\
p_{n+1} = p_n + b\left(\nabla f(x_n) - p_n\right)
\end{cases}
$$

where $a_n \in \mathbb{R}_+^d$ $b \geq 0$, $x_0, p_0 \in \mathbb{R}^d$.

# Contributions

| **Main Idea** |
|:---:|

Clipping the effective step size $a_{n+1}$ :
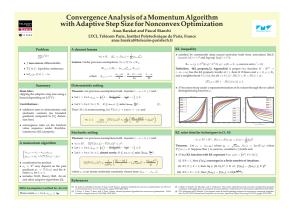
$$0 < \delta \leq a_{n+1} \leq a_{sup}(L) \tag{1}$$

| **Results** |
|:---:|

- $O(1/n)$ convergence rate for ADAM in deterministic and stochastic settings.
  (control of $\min_{0 \leq k \leq n-1} \|\nabla f(x_k)\|^2$ ).

- Convergence rate analysis of the objective function using the Kurdyka-Łojasiewicz (KŁ) property.

# Thank you for your attention

**Feel free to come to my poster**



For more details: article available on the Workshop page / arXiv.

AB, P. Bianchi. *Convergence Analysis of a Momentum Algorithm with Adaptive Step Size for Nonconvex Optimization*