

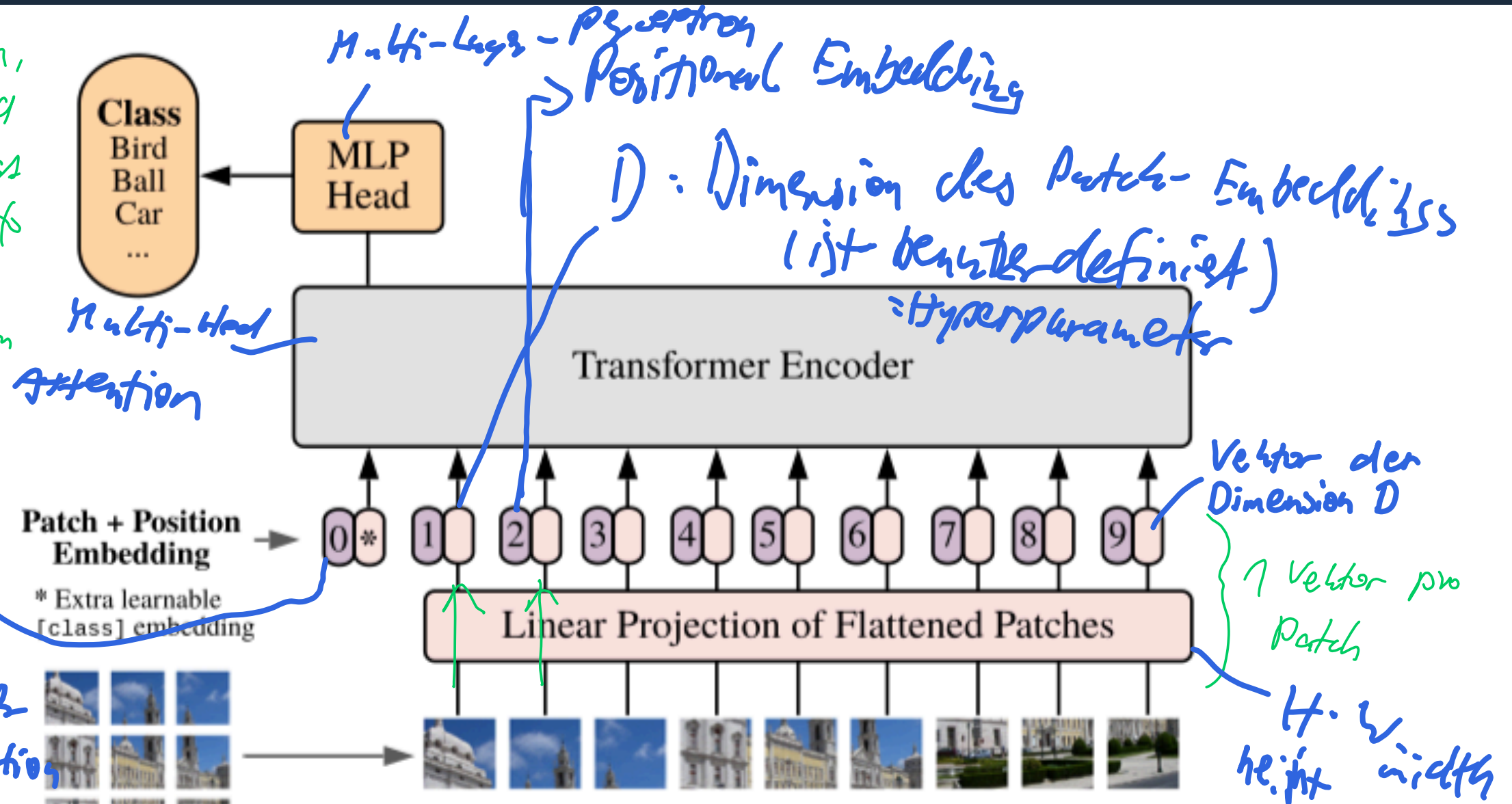
Attention und Transformer Netzwerke

hier kurz erklärt

1. Attention Konzept
- 2. Transformer Architektur**

Vision Transformer

braucht viele Daten,
wird auf 2-D Bild
auf 1-D reduziert
wird und die Info
verloren geht.
Ergebnisse trotzdem
gut



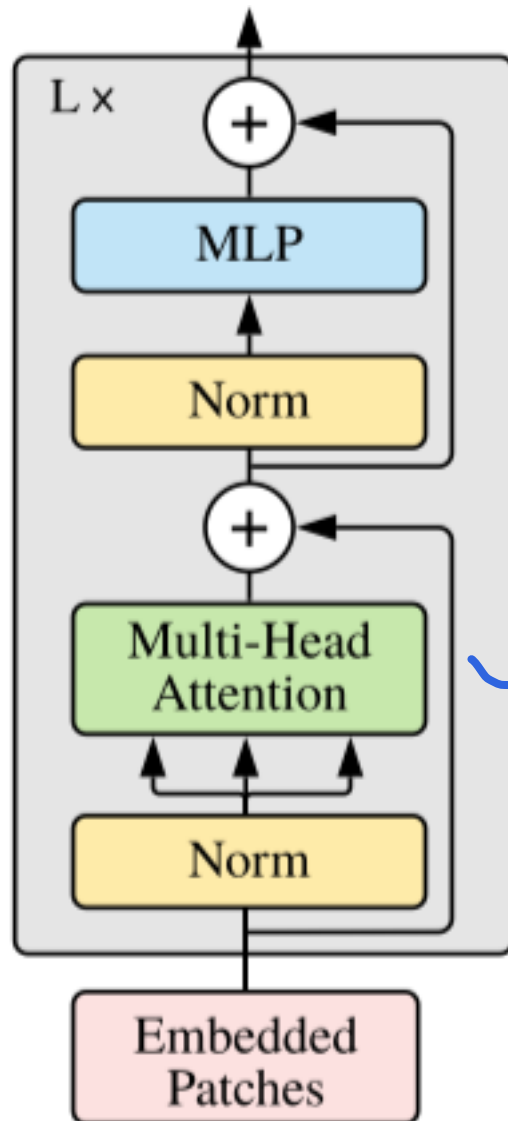
Input in
'Sequenz'
aufteilen wie
Satz bei Sprach-Tokens
L> Zerlegung des Bildes in Patches
=> Sequenzsimulation
L> Flattening => 1 Vektor

Vision Transformer

MLP

=> 7 hidden Layer
bei Pre-Training

=> 1 Layer bei Fine-Tuning
Linear Transformer Encoder



-> wiederholt sich L mal

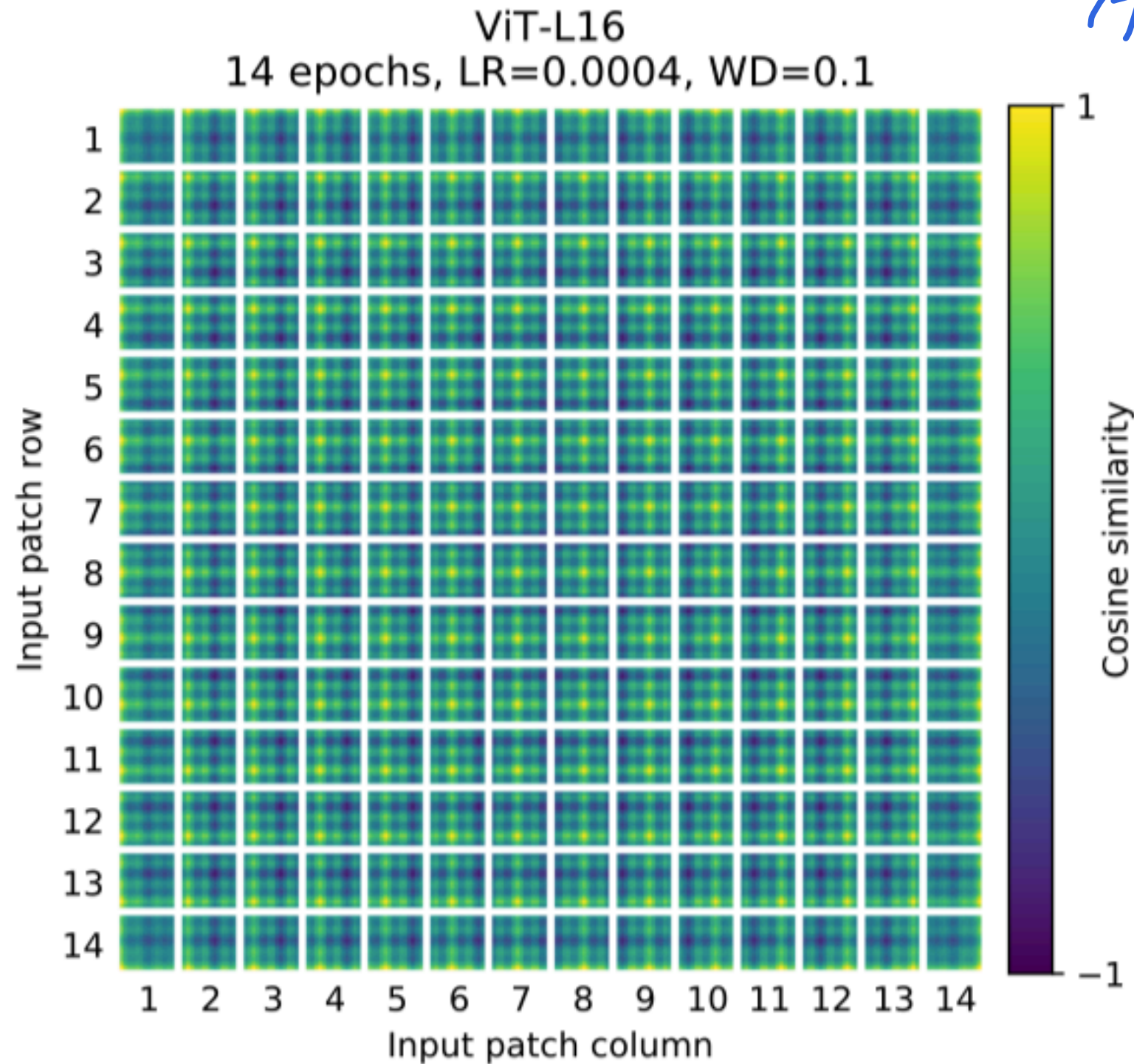
-> Residual Connection eingebaut

2 Schritte des Trainings

1) Pre-Training - nicht hochaufgelöste Bilder

2) Fine-Tuning - hochaufgelöste Bilder

Vision Transformer



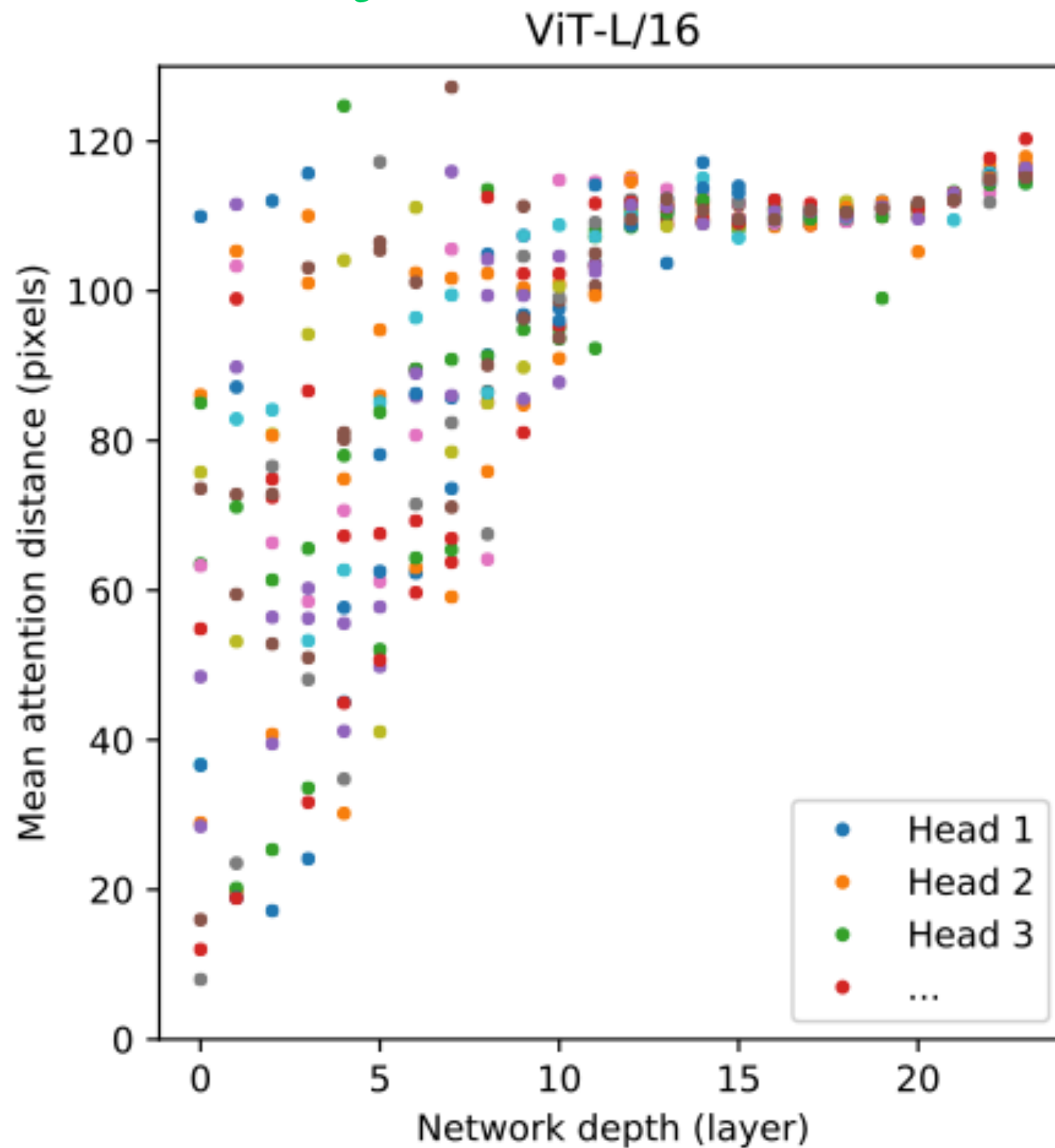
Abstände der gelernt
Positional Embeddings

=> räumlich nahe Patches
haben ähnliche Positional
Embeddings

Zusammenhang logarithmisch zwischen
räumlich weit entfernten Dingen

Vision Transformer

Nonlocal connection



→ analog to receptive fields
bei CNNs

→ große Abstände schon bei
den ersten Layern relevant

Vision Transformer



→ semantisch relevante Bildregionen
werden als wichtig markiert

