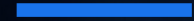


# Text 2 Video

VIDEO GENERIERUNG MIT  
BASISMODELLEN DER KI



---

# Gliederung

1. WAS IST T2V?
2. FUNKTIONSWEISE
3. ARCHITEKTURÜBERSICHT
4. TOOLS UND BEISPIELE
5. HERAUSFORDERUNGEN
6. FAZIT



---

# Was ist T2V

TEXT-TO-VIDEO(T2V) IST EINE  
TECHNOLOGIE, DIE AUS EINEM  
TEXT EIN VIDEO ERZEUGT



---

# Was bringt T2V

- REDUZIERT PRODUKTIONSZEIT & KOSTEN (FILM, WERBUNG, BILDUNG, SOCIAL MEDIA)
- VIDEOPRODUKTION OHNE KAMERA UND SCHAUSPIELER



---

# Wie funktioniert T2V?

1. TEXT DECODEN  
KI LIEST TEXT UND ERKENNT WAS  
VORKOMMT
2. NOISE GENERIEREN.  
KI STARTET MIT  
ZUFALLSVERRAUSCHTEM BILD
3. SCHRITTWEISE DENOISING  
DIFFUSSIONSMODELL ENTFERNT  
NOISE IN VIELEN KLEINEN  
SCHRITTEN
4. POSTPROCESSING  
FARBKORREKTUR, UPSCALING





---

# T2V Architekturen

- VAE (VARIATIONAL AUTOENCODER)
- GAN (GENERATIVE ADVERSARIAL NETWORK)
- DM (DIFFUSIONMODELL)
- DIT (DIFFUSION TRANSFORMER)



---

# DiT (Diffusion Transformer)

## GRUNDIDEE:

- DIFFUSION KÜMMERT SICH UM DIE BILDGENERIERUNG
- TRANSFORMER VERSTEHT ABHÄNGIGKEITEN ZWISCHEN TEXT UND ZEIT (KOHÄRENZ & PHYSIK)
- KOMBINATION ERMÖGLICHT LÄNGERE UND REALISTISCHERE VIDEOS

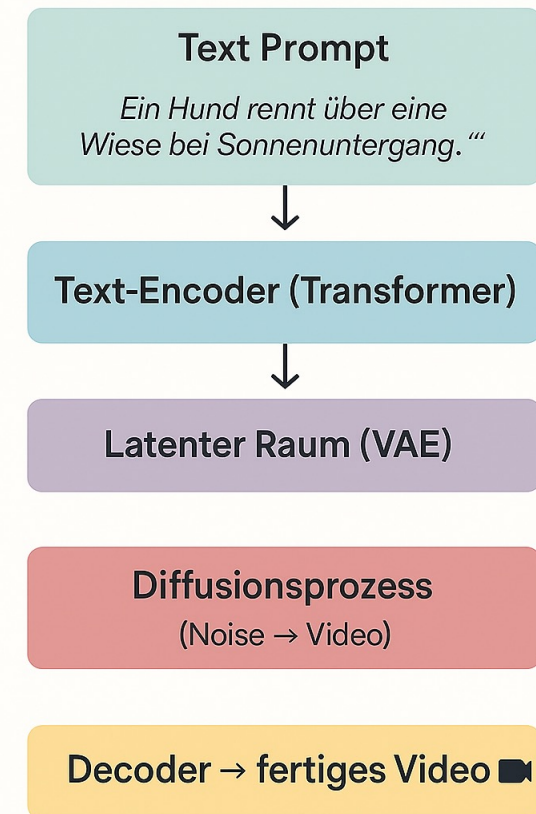


# DiT (Diffusion Transformer)

## WIE FUNKTIONIERTS ?

- ES WIRD EIN TEXT EINGEGEBEN
- DER TEXT WIRD DURCH DEN TRANSFORMER VERSTANDEN
- IM LATENTEN RAUM WIRD NACH PASSENDEN UND VORHANDENEN INFORMATIONEN GESUCHT
- DAS VIDEO WIRD DENOISED
- DAS FERTIGE VIDEO WIRD ERSTELLT

## Wie funktioniert ein Diffusion Transformer in Text-to-Video-Modell





---

# Eigenschaften von DiT

- SPATIOTEMPORAL  
UNDERSTANDING (RAUM + ZEIT)
- CROSS-ATTENTION  
MECHANISMUS (TEXT & VIDEO  
AUFEINANDER ABGEGLICHEN)
- REALISTISCHE BEWEGUNGEN
- SKALIERBARKEIT (MEHR DATEN  
+ RECHENLEISTUNG)



# Tools und Beispiele

- BEISPIELPROMPT:  
STUDENT PRÄSENTIERT EINE KI-  
PRÄSENTATION
- [SORA](#)
- [PIKA](#)
- [RUNAWAY](#)



---

# Herausforderungen

- VERSTÄNDNIS KOMPLEXER PROMPTS
- KONSISTENZ ZWISCHEN FRAMES
- PHYSIKALISCHER REALISMUS
- RECHENAUFWAND & ENERGIEVERBRAUCH
- DATEN & URHEBERRECHT
- MISSBRAUCH & DEEPPAKES



---

# Fazit

- KI KANN HEUTE AUS TEXT REALISTISCHE VIDEOS ERZEUGEN
- DIFFUSION TRANSFORMERS SIND DER AKTUELLE STAND DER TECHNIK
- NOCH GROßE HERAUSFORDERUNGEN WIE PHYSIK, ENERGIE UND ETHIK
- T2V HAT EIN GROßES POTENTIAL FÜR FILM, BILDUNG UND WERBUNG



# Quellen

- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=IV-5MZ\\_9CPY](https://www.youtube.com/watch?v=IV-5MZ_9CPY)
- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=SZ0RAJ4I-SA](https://www.youtube.com/watch?v=SZ0RAJ4I-SA)
- [HTTPS://ARXIV.ORG/PDF/2510.04999](https://arxiv.org/pdf/2510.04999)
- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=IV-5MZ\\_9CPY&T=527S](https://www.youtube.com/watch?v=IV-5MZ_9CPY&T=527S)
- CHAT GPT (STAND 10.11.2025)
- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=ZXIRUG0CN9S](https://www.youtube.com/watch?v=ZXIRUG0CN9S)
- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=DECIADFN2D4](https://www.youtube.com/watch?v=DECIADFN2D4)
- [HTTPS://WWW.YOUTUBE.COM/WATCH?V=0\\_CA0A5FMHE](https://www.youtube.com/watch?v=0_CA0A5FMHE)
- [3BLUE1BROWN\\_DL\\_PLAYLIST](#)

