



Εθνικό Μετσόβιο Πολυτεχνείο
Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών
Εργαστήριο Συστημάτων Τεχνητής Νοημοσύνης και Μάθησης

AI|LS

Παραγωγικά Μοντέλα (II)

Diffusion Models

Ηλίας Μήτσουρας

Διάρθρωση της Παρουσίασης

2

- ❑ **Εισαγωγή**
- ❑ **Diffusion Models**
 - ❑ Ιστορική Εξέλιξη
 - ❑ Διαδικασία Διάχυσης (Forward Diffusion)
 - ❑ Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)
 - ❑ Δίκτυο Αποθορυβοποίησης
 - ❑ Καθοδηγούμενη Σύνθεση Εικόνων
 - ❑ Stable Diffusion
 - ❑ Denoising Diffusion Implicit Models
 - ❑ Score-Based Generative Models
 - ❑ Noise Conditional Score Networks
 - ❑ Stochastic Differential Equations
 - ❑ Περιορισμοί
 - ❑ Generative Models Trilemma
- ❑ **Εφαρμογές**

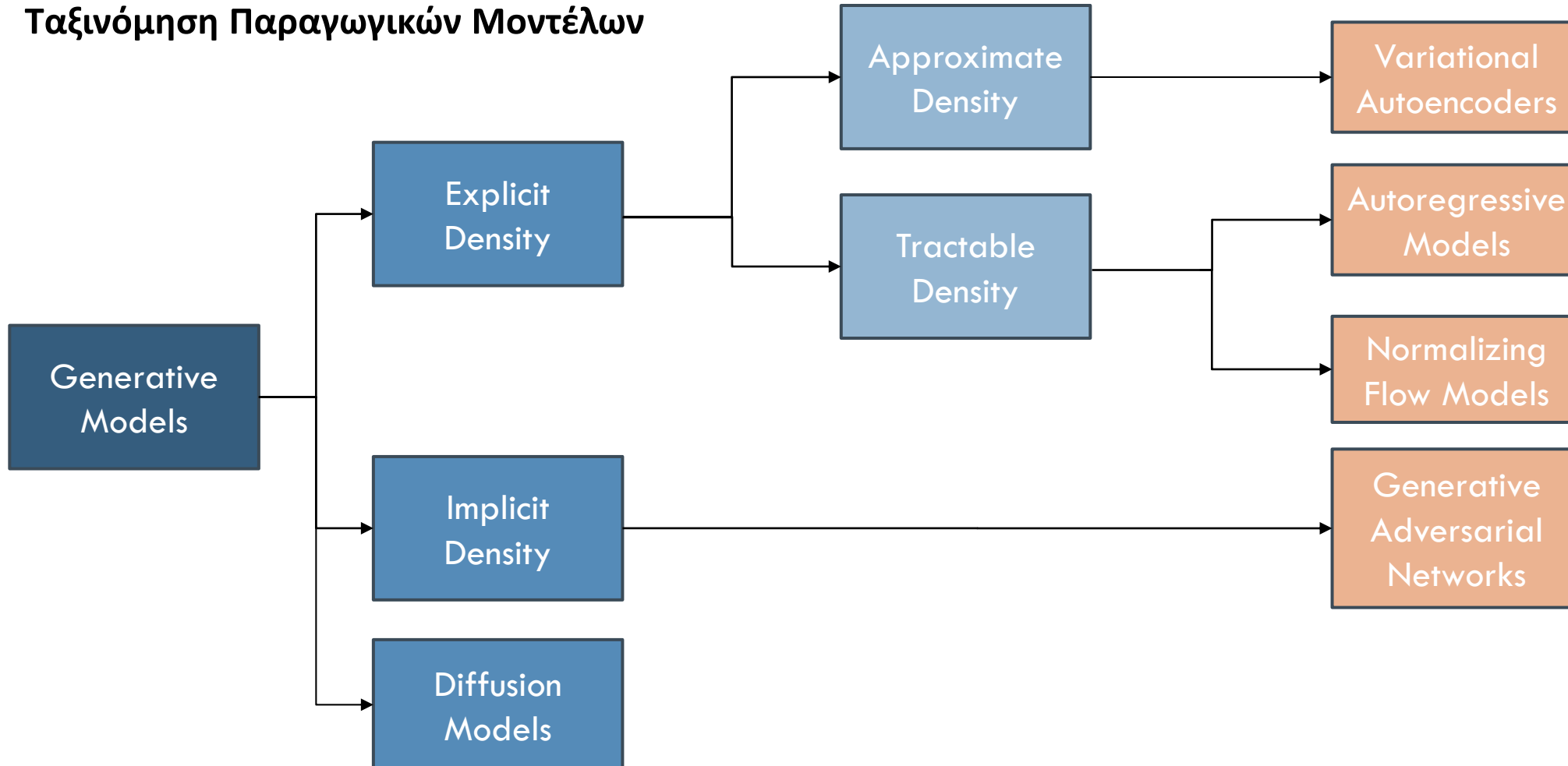
Εισαγωγή

3

Παραγωγικά Μοντέλα (*Generative Models*)

- Οικογένεια στατιστικών μοντέλων
- Σκοπός τους είναι ο προσδιορισμός των υποκείμενων μοτίβων και κατανομών των διαθέσιμων δεδομένων
→ **παραγωγή νέων δειγμάτων**, τα οποία μοιράζονται παρόμοια χαρακτηριστικά με τα αρχικά δεδομένα.
- Πιο επίσημα, δοθέντος ενός συνόλου δεδομένων $X = \{x_1, x_2, \dots, x_n\}$, και ενός συνόλου ετικετών $Y = \{y_1, y_2, \dots, y_n\}$ (προαιρετικά), ένα παραγωγικό μοντέλο προσπαθεί να προσεγγίσει:
 - την από κοινού κατανομή $p(X, Y)$ ή
 - την κατανομή $p(X)$, σε περίπτωση όπου δεν δίνεται το σύνολο Y .
- Διαφέρουν από τα διακριτικά μοντέλα (*discriminative models*), τα οποία προσπαθούν να εκτιμήσουν την υπό συνθήκη κατανομή πιθανότητας $p(Y | X)$.
- Τα παραγωγικά μοντέλα ταξινομούνται βάσει του τρόπου με τον οποίο προσπαθούν να εκτιμήσουν τη συνάρτηση πυκνότητας πιθανότητας $p_{\text{model}}(X)$.

Ταξινόμηση Παραγωγικών Μοντέλων



Διάρθρωση της Παρουσίασης

5

- ❑ Εισαγωγή
- ❑ ***Diffusion Models***
 - ❑ Ιστορική Εξέλιξη
 - ❑ Διαδικασία Διάχυσης (Forward Diffusion)
 - ❑ Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)
 - ❑ Δίκτυο Αποθρομβοποίησης
 - ❑ Καθοδηγούμενη Σύνθεση Εικόνων
 - ❑ Stable Diffusion
 - ❑ Denoising Diffusion Implicit Models
 - ❑ Score-Based Generative Models
 - ❑ Noise Conditional Score Networks
 - ❑ Stochastic Differential Equations
 - ❑ Περιορισμοί
 - ❑ Generative Models Trilemma
- ❑ Εφαρμογές

Diffusion Models

6

Ιστορική Εξέλιξη

❖ 2015

- Πρώτες ιδέες & θεμελίωση των βασικών αρχών των μοντέλων διάχυσης
- Περιορισμένη πρακτική εφαρμογή λόγω υπολογιστικών περιορισμών

❖ 2019

- Επέκταση των βασικών ιδεών και εισαγωγή της έννοιας της σταδιακής αποθορυβοποίησης

❖ 2021

- Εμφάνιση του *DALL-E* από την *OpenAI*
- Πρώτη ολοκληρωμένη προσπάθεια αντιμετώπισης του προβλήματος της παραγωγής εικόνων από κειμενικές περιγραφές

❖ 2022

- Κυκλοφορία του *Stable Diffusion* (*open source*)
- Εκρηκτική αύξηση της δημοτικότητας των μοντέλων διάχυσης

❖ 2023 – σήμερα

- Επέκταση σε νέους τομείς (βίντεο, 3D)
- Ενσωμάτωση σε εμπορικές εφαρμογές

Diffusion Models

7

Εισαγωγή

- Τα μοντέλα διάχυσης πρωτοπαρουσιάστηκαν το 2015.
- Αποτελούν μία κλάση στατιστικών παραγωγικών μοντέλων.
- Εμπνέονται από τις αρχές της θερμοδυναμικής (non-equilibrium thermodynamics).
- **Λειτουργία:** διαδοχική έγχυση θορύβου στα διαθέσιμα δεδομένα (forward diffusion) και προσπάθεια εκμάθησης της αντίστροφης διαδικασίας αποθορυβοποίησης (reverse diffusion).



Σχήμα. Διαδικασίες forward και reverse diffusion.

Diffusion Models

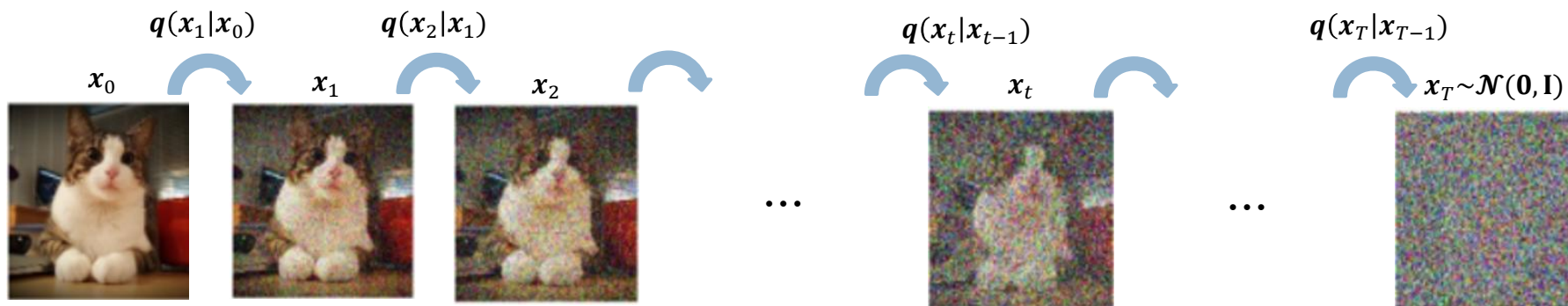
8

Διαδικασία Διάχυσης (Forward Diffusion)

- Διαδοχική προσθήκη θορύβου στα αρχικά δεδομένα μέχρι να καταστραφεί πλήρως η δομή τους \rightarrow θόρυβος.
- Η forward diffusion μοντελοποιείται μέσω μιας Μαρκοβιανής αλυσίδας.
- Δοθέντος ενός δείγματος x_0 από μια κατανομή $x_0 \sim q(x_0)$, κατά το forward diffusion προσθέτουμε διαδοχικά Γκαουσιανό θόρυβο στο δείγμα σε T βήματα, παράγοντας μια ακολουθία δειγμάτων θορύβου x_1, x_2, \dots, x_T .
- Χρησιμοποιείται ο πυρήνας μετάβασης $q(x_t|x_{t-1})$, όπου,

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I})$$

- Στην παραπάνω σχέση, το μέγεθος του βήματος ελέγχεται από μία παράμετρο, η οποία καλείται variance schedule $\{\beta_t \in (0,1)\}_{t=1}^T$.



Diffusion Models

9

Διαδικασία Διάχυσης (Forward Diffusion)

- Η από κοινού κατανομή των δειγμάτων $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T$, δεδομένου του αρχικού δείγματος \mathbf{x}_0 , δίνεται από τη σχέση,

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}),$$

η οποία προκύπτει μέσω της Μαρκοβιανής ιδιότητας.

- Στην πράξη, ο πυρήνας $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ επιλέγεται ώστε να μετασχηματίζει την αρχική κατανομή των δεδομένων σε μία εκ των προτέρων γνωστή κατανομή \rightarrow στην περίπτωσή μας κανονική κατανομή, $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.
- Αποδεικνύεται ότι:

Δειγματοληψία κατά το forward diffusion

Έστω \mathbf{x}_t ένα δείγμα το οποίο δειγματοληπτείται από την κατανομή $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ στο βήμα t . Τότε το \mathbf{x}_t μπορεί να εκφραστεί απευθείας, ως γραμμικός συνδυασμός του αρχικού δείγματος \mathbf{x}_0 και ενός διανύσματος θορύβου $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, σύμφωνα με τη σχέση,

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t,$$

όπου, $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$.

Diffusion Models

10

Διαδικασία Διάχυσης (Forward Diffusion)

- Επομένως, οι κατανομές $\mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$ και $\mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$ είναι ισοδύναμες.
- Αυτό επιτρέπει τη δειγματοληψία από οποιαδήποτε ενδιάμεση κατανομή κατά τη διάρκεια του forward diffusion, σε ένα μόνο βήμα.



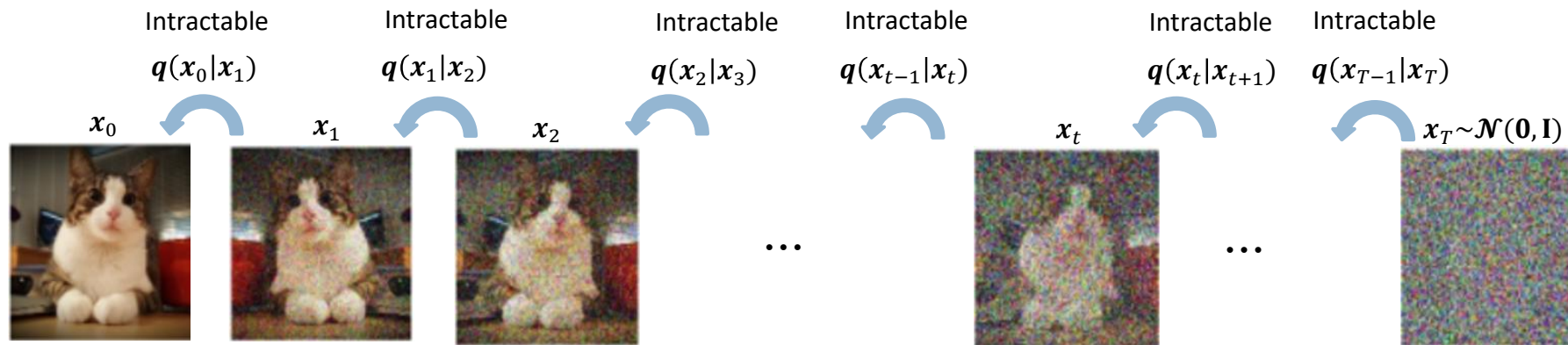
Σχήμα. Η διάχυση καταστρέφει τη δομή.

Diffusion Models

11

Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)

- Αν αντιστρέψουμε τη διαδικασία του forward diffusion και κάνουμε δειγματοληψία από την κατανομή $q(x_{t-1}|x_t)$, μπορούμε να συνθέσουμε νέα δείγματα, ξεκινώντας από Γκαουσιανό θόρυβο $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.



- Το πρόβλημα:** η κατανομή $q(x_{t-1}|x_t)$ δεν μπορεί **εκτιμηθεί** (intractable), γιατί απαιτεί υπολογισμούς με όλα τα διαθέσιμα δεδομένα. Πιο αναλυτικά,

$$q(x_{t-1}|x_t) = \frac{q(x_t|x_{t-1})q(x_{t-1})}{q(x_t)}, \text{ όπου } \boxed{q(x_t) = \int q(x_t|x_0)q(x_0)dx_0} \longrightarrow \text{Intractable}$$

Diffusion Models

12

Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)

- **Η λύση:** χρήση νευρωνικού δικτύου $p_{\theta}(x_{t-1}|x_t)$ για την προσέγγιση των υπό συνθήκη κατανομών $q(x_{t-1}|x_t)$, δηλαδή,

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t)),$$

όπου η παράμετρος θ υποδηλώνει τις παραμέτρους του δικτύου.

- **Πρόβλημα:** οι ποσότητες $\mu_{\theta}(x_t, t)$ και $\Sigma_{\theta}(x_t, t)$ δεν είναι γνωστές, ώστε να χρησιμοποιηθούν ως ground truth για την εκπαίδευση του δικτύου.
- **Υπενθύμιση:** η κατανομή $q(x_{t-1}|x_t)$ δεν μπορεί **εκτιμηθεί** (intractable).
- Αποδεικνύεται, ότι εάν δεσμεύσουμε την κατανομή αυτή στο δείγμα x_0 , αυτή μπορεί πλέον να εκτιμηθεί. Δηλαδή,

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}_t(x_t, x_0), \tilde{\beta}_t \mathbf{I}).$$

- Εφαρμόζοντας τον κανόνα του Bayes για την κατανομή $q(x_{t-1}|x_t, x_0)$, λαμβάνουμε έπειτα από υπολογισμούς:

$$\tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} x_0 \text{ και } \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t.$$

Diffusion Models

13

Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)

- Αν τώρα μετασχηματίσουμε την εξίσωση που χρησιμοποιούμε για τη δειγματοληψία κατά το forward diffusion, λαμβάνουμε,

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t \Rightarrow \mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\varepsilon}_t}{\sqrt{\bar{\alpha}_t}}.$$

- Αντικαθιστώντας την ποσότητα αυτή στη σχέση που υπολογίσαμε για τη μέση τιμή $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$ προκύπτει ότι,

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_t \right).$$

- Επομένως, αρκεί να εκπαιδεύσουμε το δίκτυό μας ώστε να προσεγγίζει τη μέση τιμή $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$.
- Ισοδύναμα, αρκεί να εκπαιδεύσουμε το δίκτυό μας ώστε να προσεγγίζει το θόρυβο $\boldsymbol{\varepsilon}_t$ για την εκάστοτε είσοδο \mathbf{x}_t , στο τυχόν βήμα t . Έπειτα, η μέση τιμή $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$ θα υπολογίζεται προσεγγιστικά μέσω της σχέσης,

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) \right),$$

όπου $\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$ η εκτίμηση του θορύβου $\boldsymbol{\varepsilon}_t$ στο βήμα t .

Diffusion Models

14

Δίκτυο Αποθρομβοποίησης

- Οι Ho et al. όρισαν τη συνάρτηση απώλειας L_t του δικτύου αποθρομβοποίησης, ώστε να ελαχιστοποιείται η διαφορά μεταξύ της εκτίμησης $\mu_\theta(\mathbf{x}_t, t)$ και της μέσης τιμής $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0)$, σύμφωνα με τη σχέση,

$$\begin{aligned} L_t &= \mathbb{E}_{\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{1}{2 \|\Sigma_\theta(\mathbf{x}_t, t)\|_2^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] \\ &= \mathbb{E}_{\mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\frac{(1 - a_t)^2}{2 a_t (1 - \bar{a}_t) \|\Sigma_\theta(\mathbf{x}_t, t)\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{a}_t} \mathbf{x}_0 + \sqrt{1 - \bar{a}_t} \epsilon_t, t)\|^2 \right]. \end{aligned}$$

- Παρόλα αυτά, διαπίστωσαν ότι παραλείποντας τον πρώτο όρο στην παραπάνω σχέση, απλοποιείται σημαντικά η συνάρτηση απώλειας και βελτιώνεται η ποιότητα των παραγόμενων δειγμάτων:

$$L_t^{simple} = \mathbb{E}_{t \sim \mathcal{U}[1, T], \mathbf{x}_0 \sim \mathbf{q}(\mathbf{x}_0), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\|\epsilon_t - \epsilon_\theta(\sqrt{\bar{a}_t} \mathbf{x}_0 + \sqrt{1 - \bar{a}_t} \epsilon_t, t)\|^2 \right].$$

Diffusion Models

15

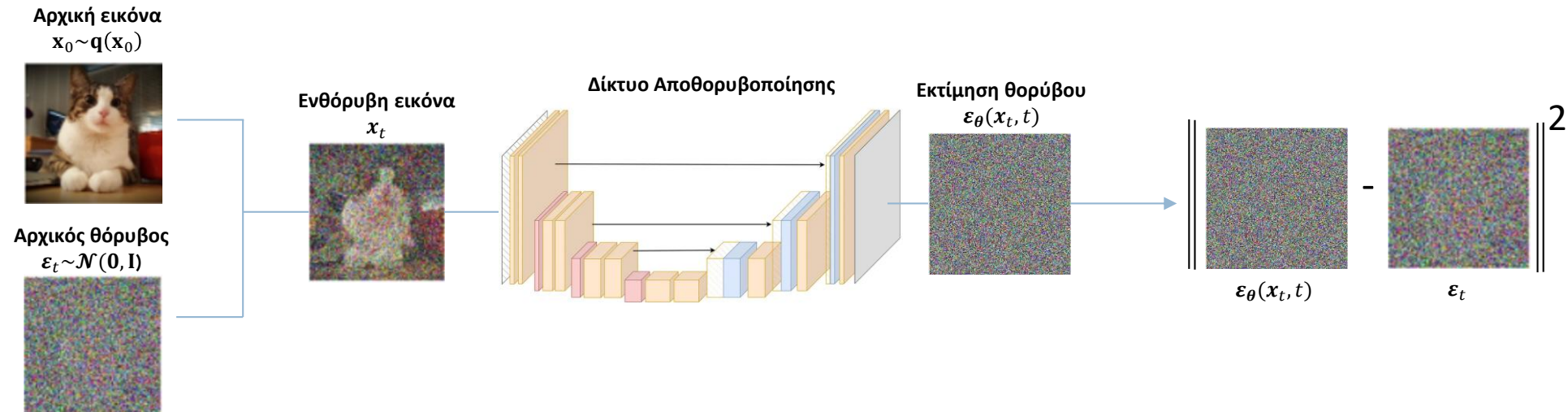
Δίκτυο Αποθουροποίησης

■ Διαδικασία εκπαίδευσης

Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
      $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$ 
6: until converged
  
```



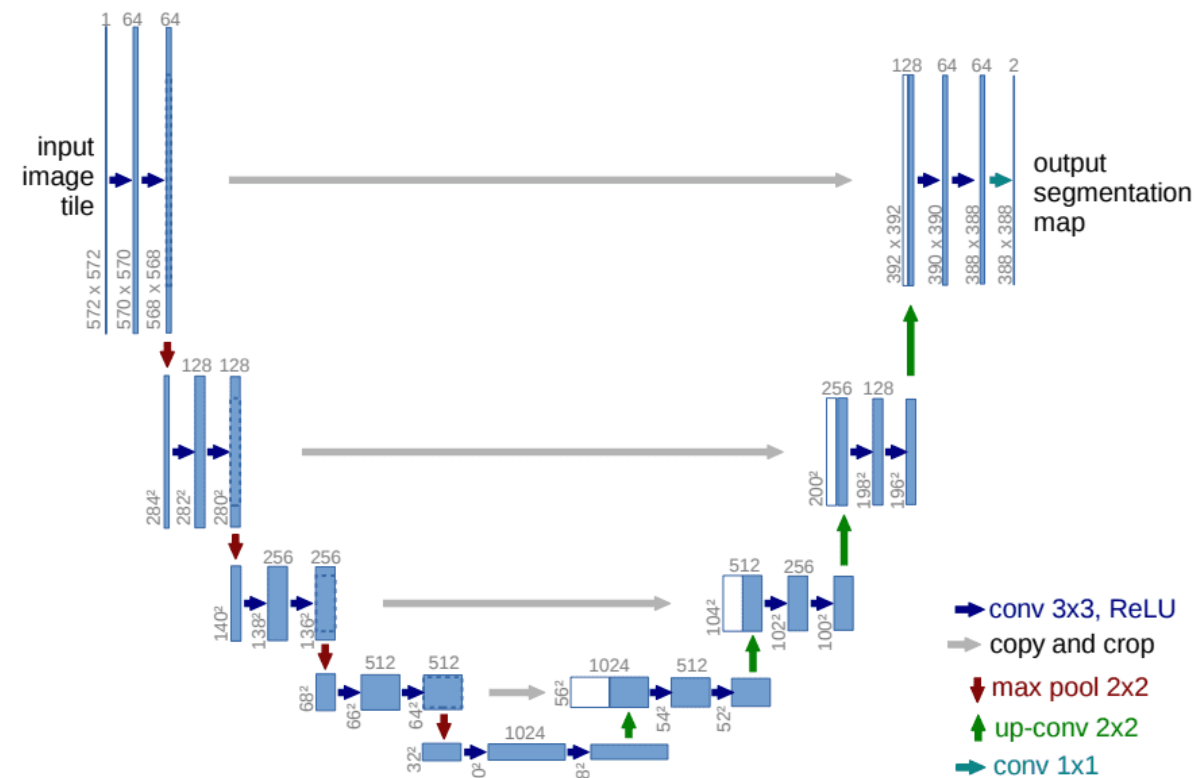
Diffusion Models

16

Δίκτυο Αποθρυβοποίησης

□ Αρχιτεκτονική

- Η πιο συνηθισμένη αρχιτεκτονική του δικτύου αποθρυβοποίησης είναι αυτή του δικτύου U-Net, η οποία περιλαμβάνει εκτός από συνελκτικά layers και attention layers.
- Οι πιο καινούριες αρχιτεκτονικές βασίζονται εξ ολοκλήρου σε δομές transformers, οι οποίες είναι ειδικά σχεδιασμένες για τις ανάγκες των diffusion μοντέλων (Diffusion transformers).



Diffusion Models

17

Δίκτυο Αποθουβοποίησης

■ Σύνθεση νέων δειγμάτων

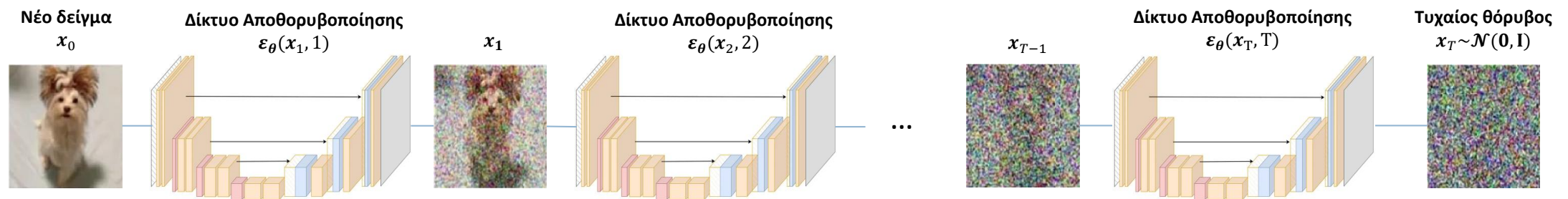
- Αξιοποιώντας το εκπαιδευμένο δίκτυο αποθουβοποίησης $\epsilon_{\theta}(\mathbf{x}_t, t)$, μπορούμε να ακολουθήσουμε τη reverse diffusion διαδικασία βήμα προς βήμα και να παράξουμε νέα δεδομένα.

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 

```



Diffusion Models

18

Καθοδηγούμενη Σύνθεση Εικόνων

- Καθοδήγηση της διαδικασίας σύνθεσης των εικόνων βάσει εξωτερικών συνθηκών (π.χ. κειμενικές περιγραφές).
- Ας υποθέσουμε ότι θέλουμε να παράξουμε νέα δείγματα βάσει της συνθήκης $\mathbf{y} \rightarrow$ ετικέτα της κλάσης της εικόνας.
- Οι Sohl-Dickstein et al. και αργότερα οι Dhariwal και Nichol έδειξαν ότι μπορούμε να χρησιμοποιήσουμε ένα μοντέλο ταξινόμησης $\mathbf{p}_\phi(\mathbf{y}|\mathbf{x}_t, t)$, για να καθοδηγήσουμε τη διαδικασία του reverse diffusion προς την επιθυμητή κλάση \mathbf{y} .
- Για να το επιτύχουμε αυτό, εκπαιδεύουμε τον ταξινομητή $\mathbf{p}_\phi(\mathbf{y}|\mathbf{x}_t, t)$ στις θορυβώδεις εικόνες \mathbf{x}_t , ώστε να προβλέπει τις κλάσεις τους \mathbf{y} .
- Έπειτα, αντί να λαμβάνουμε δείγματα από την κατανομή $\mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$, λαμβάνουμε δείγματα από την κατανομή $\mathcal{N}(\hat{\boldsymbol{\mu}}(\mathbf{x}_t|\mathbf{y}, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$, όπου,

$$\hat{\boldsymbol{\mu}}(\mathbf{x}_t|\mathbf{y}, t) = \boldsymbol{\mu}_\theta(\mathbf{x}_t, t) + s\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)\nabla_{\mathbf{x}_t}\log \mathbf{p}_\phi(\mathbf{y}|\mathbf{x}_t, t).$$

- Η μέθοδος αυτή καλείται **Classifier Guidance**.
- Αν και αποτελεσματική, απαιτεί την εκπαίδευση ενός επιπλέον μοντέλου ταξινόμησης στο ίδιο σύνολο δεδομένων \rightarrow Πολύ υψηλό υπολογιστικό κόστος.

Diffusion Models

19

Καθοδηγούμενη Σύνθεση Εικόνων

- Για να αντιμετωπίσουν το πρόβλημα αυτό, οι Ho et al. απέδειξαν ότι η καθοδήγηση μέσω της συνθήκης \mathbf{y} μπορεί να υλοποιηθεί και χωρίς τη χρήση του εξωτερικού ταξινομητή $\mathbf{p}_\varphi(\mathbf{y}|\mathbf{x}_t, t)$.
- **Λύση:** από κοινού εκπαίδευση ενός **υπό συνθήκη** δικτύου αποθρομβοποίησης $\varepsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)$ και ενός **απλού** δικτύου αποθρομβοποίησης $\varepsilon_\theta(\mathbf{x}_t, \mathbf{y} = \emptyset, t)$. Στην πραγματικότητα πρόκειται για **το ίδιο δίκτυο**.
- Η εκπαίδευση πραγματοποιείται με μικτό τρόπο \rightarrow τυχαία επιλογή μεταξύ ετικετών \mathbf{y} και $\mathbf{y} = \emptyset$, ούτως ώστε το μοντέλο να εκτεθεί τόσο στην υπό συνθήκη όσο και στην χωρίς συνθήκη παραγωγή.
- Για τη παραγωγή νέων δειγμάτων ακολουθείται και πάλι η ίδια διαδικασία (Αλγόριθμος 2), με τη διαφορά ότι:

$$\tilde{\varepsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t) = \varepsilon_\theta(\mathbf{x}_t, \mathbf{y} = \emptyset, t) + w(\varepsilon_\theta(\mathbf{x}_t, \mathbf{y}, t) - \varepsilon_\theta(\mathbf{x}_t, \mathbf{y} = \emptyset, t)),$$

όπου w μία παράμετρος η οποία καλείται **guidance scale** και η οποία ορίζει το βαθμό της επίδρασης της συνθήκης \mathbf{y} στη διαδικασία της σύνθεσης.

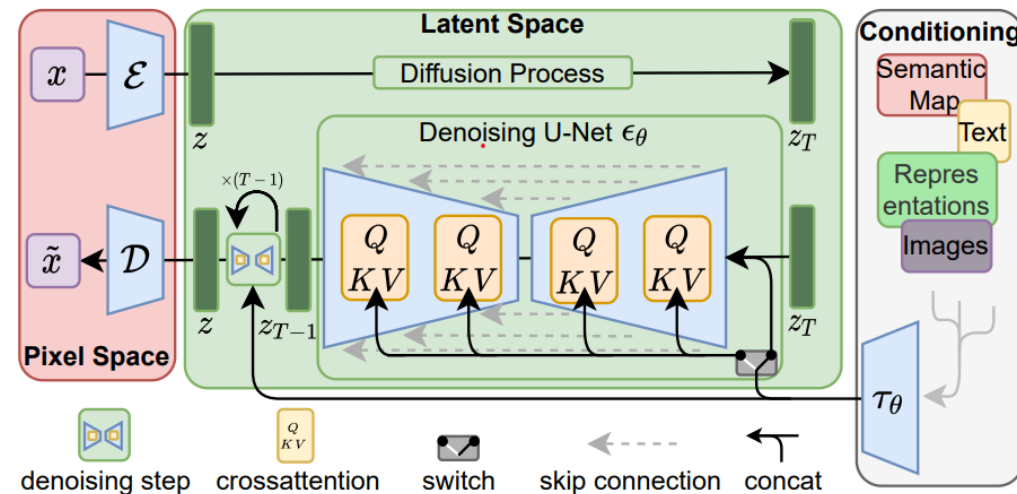
- **Ερώτημα:** πώς η συνθήκη \mathbf{y} εισέρχεται στο δίκτυο αποθρομβοποίησης ε_θ ;

Diffusion Models

20

Stable Diffusion

- Πρόκειται για ένα latent μοντέλο διάχυσης \rightarrow διαδικασία διάχυσης και αντίστροφης διάχυσης σε latent χώρο.
- Χρησιμοποιεί ένα δίκτυο αποθρυβοποίησης $\epsilon_{\theta}(\mathbf{z}_t, \mathbf{y}, t)$ τύπου U-Net, όπου $\mathbf{z}_t = \mathcal{E}(\mathbf{x})$ και $\mathbf{x} \sim \mathbf{q}(\mathbf{x}_0)$.
- Κάθε block του U-Net αποτελείται από συνελκτικά layers ακολουθούμενα από self-attention και cross-attention layers.
- Η κειμενική καθοδήγηση \mathbf{y} εισέρχεται στη διαδικασία της διάχυσης μέσω των cross-attention layers.



Diffusion Models

21

Denoising Diffusion Implicit Models

- Στα παραδοσιακά μοντέλα διάχυσης (DDPMs) η διαδικασία του forward diffusion μπορεί να περιλαμβάνει αρκετές χιλιάδες βήματα.
- Για να παράξουμε νέα δείγματα, θα πρέπει να ακολουθήσουμε αντίστροφα όλα τα βήματα αυτά → υψηλοί χρόνοι εκτέλεσης.
- **Λύση:** τα Denoising Diffusion Implicit Models, τα οποία επεκτείνουν την ιδέα των παραδοσιακών μοντέλων διάχυσης σε μη-Μαρκοβιανές περιπτώσεις.

Εκπαίδευση	Δειγματοληψία
Μεγάλος αριθμός χρονικών βημάτων T (συνήθως μερικές χιλιάδες)	Υποσύνολο των βημάτων που χρησιμοποιούνται κατά την εκπαίδευση

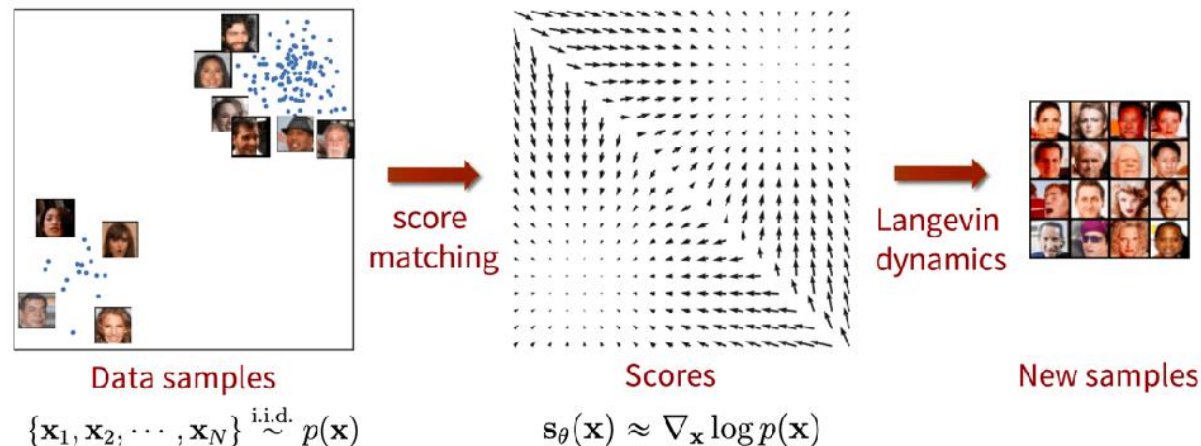
Diffusion Models

22

Score-Based Generative Models

- Η βασική ιδέα των μοντέλων αυτών είναι η χρήση ενός νευρωνικού δικτύου s_θ για την εκτίμηση της **score function** και η δημιουργία νέων δειγμάτων με χρήση score-based μεθόδων δειγματοληψίας.
- Η score function μιας πυκνότητας πιθανότητας $p(x)$ ορίζεται ως $\nabla_x \log p(x)$.
- Πρόκειται ουσιαστικά για ένα διανυσματικό πεδίο, το οποίο δείχνει προς την κατεύθυνση εκείνη στην οποία μεγιστοποιείται η λογαριθμική πυκνότητα των δεδομένων.
- Το νευρωνικό δίκτυο εκπαιδεύεται ώστε να προσεγγίζει την ποσότητα αυτή, δηλαδή $s_\theta(x) \approx \nabla_x \log p(x)$, σύμφωνα με τη συνάρτηση απώλειας,

$$\mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2]$$



Diffusion Models

23

Score-Based Generative Models

■ Δειγματοληψία με Langevin Dynamics

- Αφού έχουμε εκπαιδεύσει το δίκτυο s_θ , μπορούμε να παράξουμε νέα δείγματα χρησιμοποιώντας μια Markov Chain Monte Carlo μέθοδο δειγματοληψίας, τα **Langevin Dynamics**.
- Η μέθοδος αυτή, επιτρέπει τη λήψη δειγμάτων από μία κατανομή $p(x)$, χρησιμοποιώντας αποκλειστικά τη score function $\nabla_x \log p(x)$.
- Πιο συγκεκριμένα, ξεκινώντας από ένα δείγμα $x_0 \sim \pi(x)$, όπου $\pi(x)$ μια prior κατανομή και με βήμα $\epsilon > 0$, η μέθοδος των Langevin Dynamics υπολογίζει αναδρομικά, για T βήματα, το ακόλουθο:

$$\tilde{x}_t = \tilde{x}_{t-1} + \frac{\epsilon}{2} \nabla_x \log p(\tilde{x}_{t-1}) + \sqrt{\epsilon} z_t$$

όπου $z_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

- Όταν $\epsilon \rightarrow 0$ και $T \rightarrow \infty$, η κατανομή των δειγμάτων \tilde{x}_T προσεγγίζει την κατανομή $p(x)$.
- Στην πράξη, επιλέγουμε ένα αρκετά μικρό ϵ και ένα αρκετά μεγάλο T .

Diffusion Models

24

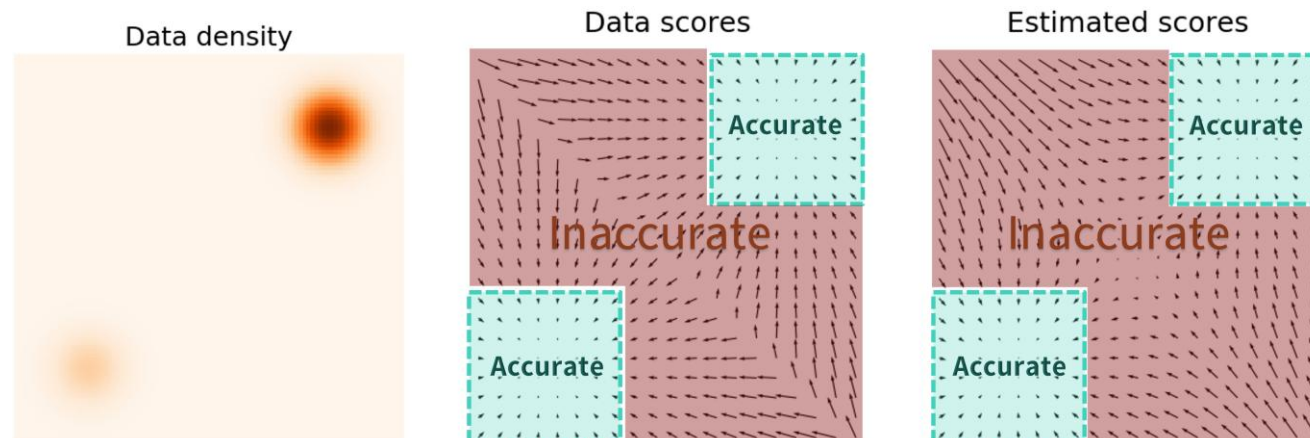
Score-Based Generative Models

■ Περιοχές Χαμηλής Πυκνότητας Δεδομένων (Low Data Density Regions)

- Το βασικό πρόβλημα των score-based generative μοντέλων είναι ότι η προσέγγιση της score function δεν είναι ακριβής σε περιοχές χαμηλής συγκέντρωσης των δεδομένων.
- Τα score-based generative μοντέλα εκπαιδεύονται βάσει της συνάρτησης απώλειας,

$$\mathbb{E}_{p(x)} [\|\nabla_x \log p(x) - s_\theta(x)\|_2^2] = \int p(x) \|\nabla_x \log p(x) - s_\theta(x)\|_2^2 dx.$$

- Επομένως, οι περιοχές χαμηλής συγκέντρωσης αμελούνται κατά την εκπαίδευση, λόγω του όρου $p(x)$.

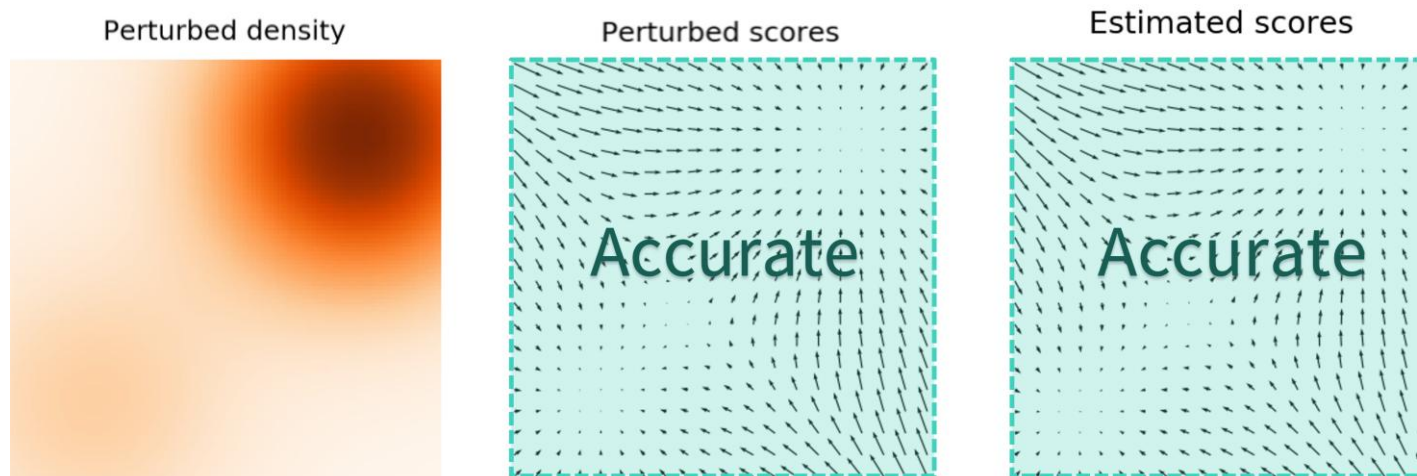


Diffusion Models

25

Noise Conditional Score Networks

- Αποτελούν είδος score-based generative μοντέλων, τα οποία αναπτύχθηκαν για να αντιμετωπιστεί το πρόβλημα των low data density regions.
- Χρησιμοποιούν ένα νευρωνικό δίκτυο για την εκτίμηση της score function της κατανομής των δεδομένων **έπειτα από την προσθήκη θορύβου**.
- Με την προσθήκη του θορύβου ομαλοποιείται η πυκνότητα των περιοχών και αποφεύγονται οι περιοχές χαμηλής πυκνότητας.



Diffusion Models

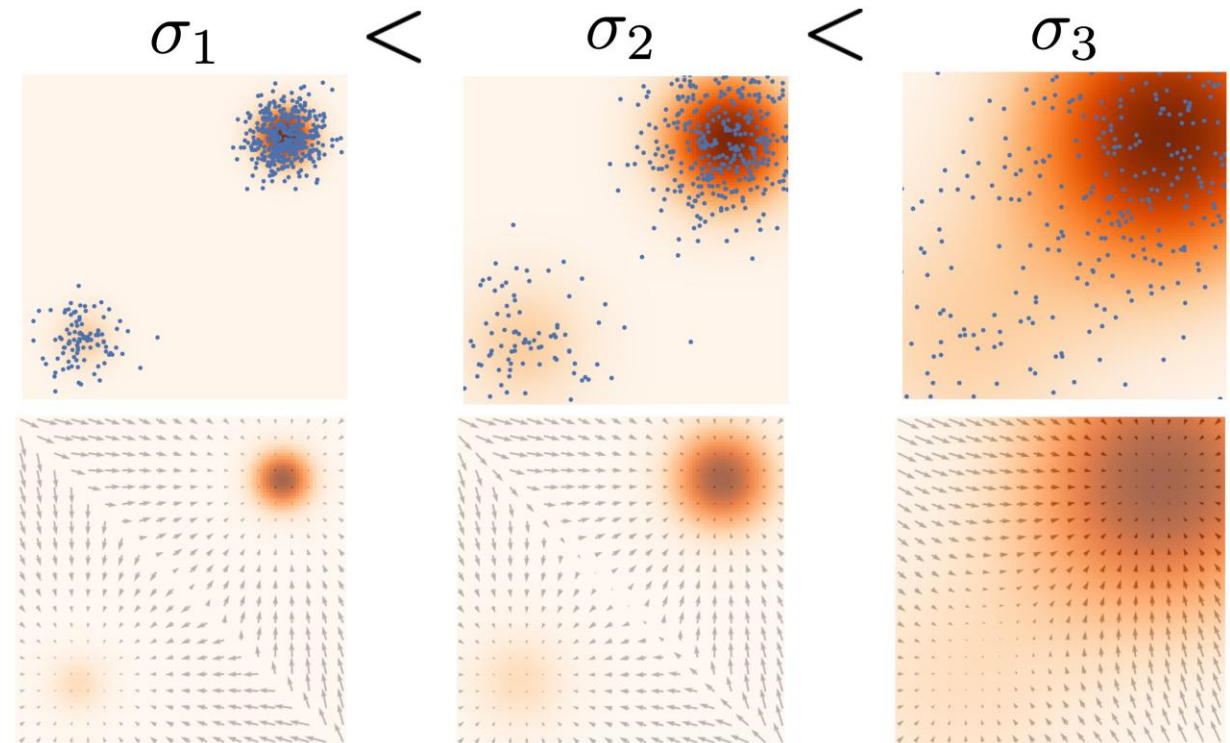
26

Noise Conditional Score Networks

- Έστω $\sigma_1 < \sigma_2 < \dots < \sigma_L$ μια ακολουθία επιπέδων θορύβου.
- Προσθέτουμε Γκαουσιανό θόρυβο $\mathcal{N}(\mathbf{0}, \sigma_i^2 \mathbf{I}), i = 1, 2, \dots, L$ στα δεδομένα και λαμβάνουμε την ενθόρυβη κατανομή,

$$\begin{aligned} p_{\sigma_i}(\mathbf{x}) &= \int p(\mathbf{y}) p(\mathbf{x}|\mathbf{y}) d\mathbf{y} \\ &= \int p(\mathbf{y}) \mathcal{N}(\mathbf{x}; \mathbf{y}, \sigma_i^2 \mathbf{I}) d\mathbf{y}. \end{aligned}$$

- Εκπαιδεύουμε ένα Noise Conditional score-based νευρωνικό δίκτυο $\mathbf{s}_{\theta}(\mathbf{x}, \sigma_i)$, το προσεγγίζει αυτή τη score function, δηλαδή $\mathbf{s}_{\theta}(\mathbf{x}, \sigma_i) \approx \nabla_{\mathbf{x}} \log p_{\sigma_i}(\mathbf{x})$, για κάθε $i = 1, 2, \dots, L$.



Diffusion Models

27

Noise Conditional Score Networks

- Σε κάθε επίπεδο θορύβου, η αντικειμενική συνάρτηση ισούται με,

$$\ell(\boldsymbol{\theta}; \sigma) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 \mathbf{I})} \left[\left\| \mathbf{s}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}, \sigma) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|_2^2 \right]$$

- Συνδυάζοντας την παραπάνω συνάρτηση για κάθε $\sigma \in \{\sigma_i\}_{i=1}^L$, λαμβάνουμε την ολική αντικειμενική συνάρτηση εκπαίδευσης του δικτύου $\mathbf{s}_{\boldsymbol{\theta}}(\mathbf{x}, \sigma)$:

$$\mathcal{L}(\boldsymbol{\theta}, \{\sigma_i\}_{i=1}^L) = \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) \ell(\boldsymbol{\theta}; \sigma_i),$$

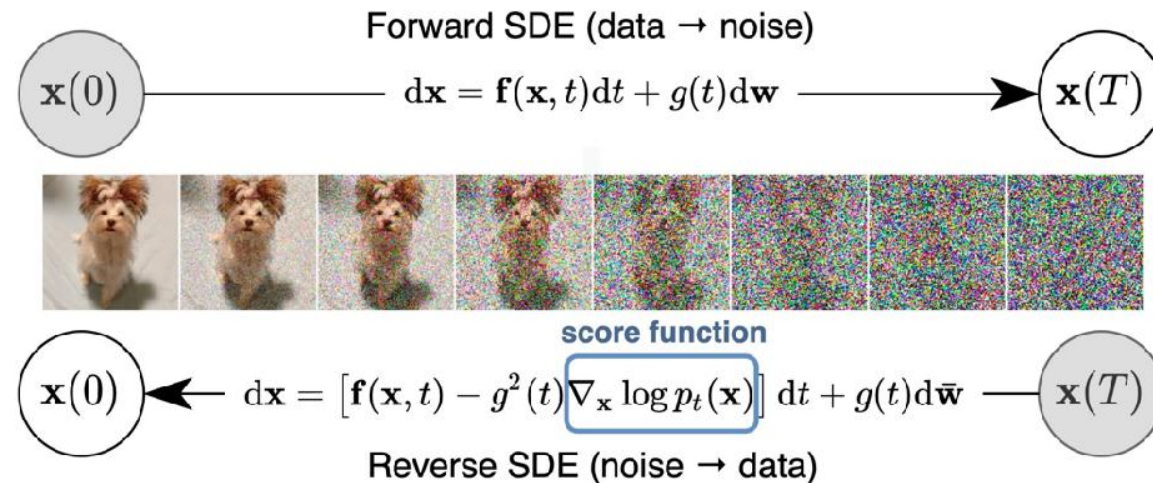
όπου $\lambda(\cdot)$ είναι ένας συντελεστής, ο οποίος συχνά επιλέγεται ίσος με $\lambda(\sigma) = \sigma^2$.

Diffusion Models

28

Stochastic Differential Equations

- Τόσο τα παραδοσιακά μοντέλα διάχυσης, όσο και τα score-based generative μοντέλα, μπορούν να θεωρηθούν ως διακριτοποιήσεις στοχαστικών διαφορικών εξισώσεων.
- Στις έως τώρα προσεγγίσεις, χρησιμοποιείται ένα πεπερασμένο πλήθος επιπέδων θορύβου.
- Μπορούμε να επεκτείνουμε τη λογική αυτή και να χρησιμοποιήσουμε ένα συνεχές φάσμα κατανομών θορύβου για να προσθέσουμε θόρυβο στα δεδομένα.
- Οι κατανομές των ενθόρυβων δεδομένων, εξελίσσονται χρονικά σύμφωνα με μία Στοχαστική Διαφορική Εξίσωση:



Diffusion Models

29

Stochastic Differential Equations

- Παρατηρούμε ότι υπάρχουν αρκετές αναλογίες μεταξύ των παραδοσιακών μοντέλων διάχυσης, των Score-based Generative μοντέλων και των SDEs.
 - Τα διακριτά βήματα $t = 1, 2, \dots, T$, γενικεύονται σε βήματα στο συνεχές φάσμα $t \in [0, T]$.
 - Τα διακριτά επίπεδα θορύβου που ελέγχονται μέσω της παραμέτρου β_t , αντιστοιχούν στο συντελεστή διάχυσης $g(t)$ της SDE.
- Προκειμένου τώρα να παράξουμε νέα δείγματα $\mathbf{x}(\mathbf{0}) \sim \mathbf{p}_0$, μπορούμε, ξεκινώντας από ένα δείγμα $\mathbf{x}(T) \sim \mathbf{p}_T$ να ακολουθήσουμε την αντίστροφη SDE,

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t) - \mathbf{g}(t)^2 \nabla_{\mathbf{x}} \log_t \mathbf{p}(\mathbf{x})]dt + \mathbf{g}(t)d\bar{\mathbf{w}}.$$

- Για να κάνουμε δειγματοληψία από την κατανομή \mathbf{p}_0 , θα πρέπει να προσομοιώσουμε την αντίστροφη SDE, η οποία προϋποθέτει την εκτίμηση της score function $\nabla_{\mathbf{x}} \log_t \mathbf{p}(\mathbf{x})$.
- Η score function αυτή είναι προφανώς η ίδια με αυτή της περίπτωσης των Score-based Generative Μοντέλων.

Diffusion Models

30

Περιορισμοί

- ❑ Υψηλό Υπολογιστικό Κόστος
 - Ανάγκη για πληθώρα υπολογιστικών πόρων → δύσκολη υλοποίηση σε πραγματικό χρόνο
 - Αργή και ιδιαίτερα απαιτητική σε hardware εκπαίδευση
- ❑ Δυσκολία Γενίκευσης σε Άγνωστα Δεδομένα
 - Ανάγκη επανεκπαίδευσης και προσαρμογής
- ❑ Υψηλό Επίπεδο Πολυπλοκότητας
 - Αδυναμία ενσωμάτωσης και χρήσης σε εφαρμογές οι οποίες απαιτούν γνώση του ακριβούς τρόπου λειτουργίας τους.
- ❑ Ελλιπής ικανότητα ερμηνείας των αποτελεσμάτων
 - Προκλήσεις σε εφαρμογές ερμηνεύσιμης τεχνητής νοημοσύνης (Explainable AI)
- ❑ Πολωμένα Αποτελέσματα
 - Οι παραγόμενες εικόνες μπορεί να παρουσιάζουν απόκλιση από τις κειμενικές περιγραφές, λόγω έλλειψης ισορροπίας στα δεδομένα εκπαίδευσης.

"A group of watches showing 5 minutes past 12"



Παραγωγή με χρήση του Stable Diffusion

"A left-handed person writing down on a notebook"



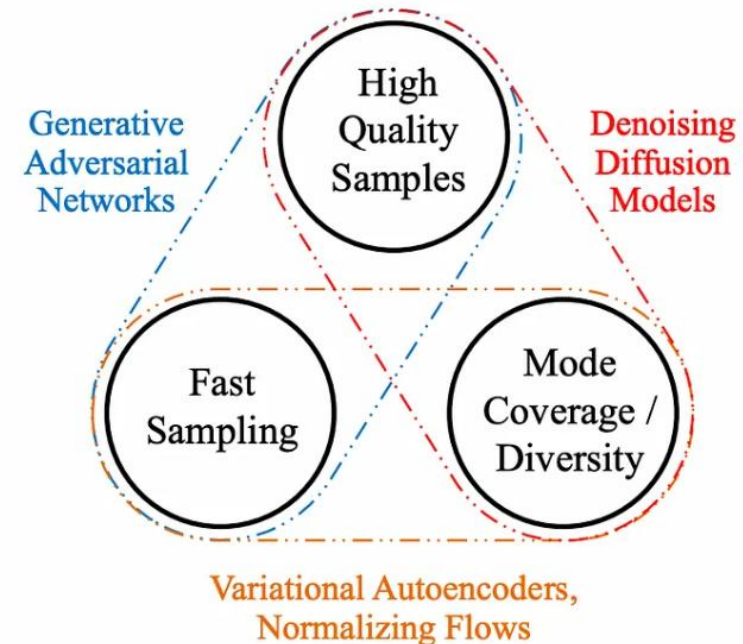
Παραγωγή με χρήση του Stable Diffusion

Diffusion Models

31

Generative Models Trilemma

- Βασικοί στόχοι ενός παραγωγικού μοντέλου:
 - 1) Παραγωγή υψηλής ποιότητας δειγμάτων,
 - 2) Παραγωγή diverse δειγμάτων και
 - 3) Γρήγορη παραγωγή δειγμάτων με χαμηλό υπολογιστικό κόστος.
- Δεν υπάρχει κάποιο παραγωγικό framework που να ικανοποιεί και τα 3 κριτήρια.
- Variational Autoencoders (VAEs) → χαμηλή ποιότητα παραγόμενων δειγμάτων.
- Generative Adversarial Networks (GANs) → περιορισμένη ποικιλία στα παραγόμενα δείγματα.
- Diffusion Models → υψηλό υπολογιστικό κόστος (μεγάλος αριθμός βημάτων για παραγωγή νέων δειγμάτων).



Διάρθρωση της Παρουσίασης

32

- ❑ Εισαγωγή
- ❑ Diffusion Models
 - ❑ Ιστορική Εξέλιξη
 - ❑ Διαδικασία Διάχυσης (Forward Diffusion)
 - ❑ Διαδικασία Αντίστροφης Διάχυσης (Reverse Diffusion)
 - ❑ Δίκτυο Αποθορυβοποίησης
 - ❑ Καθοδηγούμενη Σύνθεση Εικόνων
 - ❑ Stable Diffusion
 - ❑ Denoising Diffusion Implicit Models
 - ❑ Score-Based Generative Models
 - ❑ Noise Conditional Score Networks
 - ❑ Stochastic Differential Equations
 - ❑ Περιορισμοί
 - ❑ Generative Models Trilemma
- ❑ **Εφαρμογές**

Text-to-Image Synthesis

- Τα μοντέλα διάχυσης μπορούν να χρησιμοποιηθούν για την **παραγωγή εικόνων** από κειμενικές περιγραφές.
- Γνωστά μοντέλα:
 - *DALL-E 2, DALL-E 3*
 - *Stable Diffusion*
 - *MidJourney*
 - *Imagen*



"a hedgehog using a calculator"



"a corgi wearing a red bowtie and a purple party hat"



"robots meditating in a vipassana retreat"



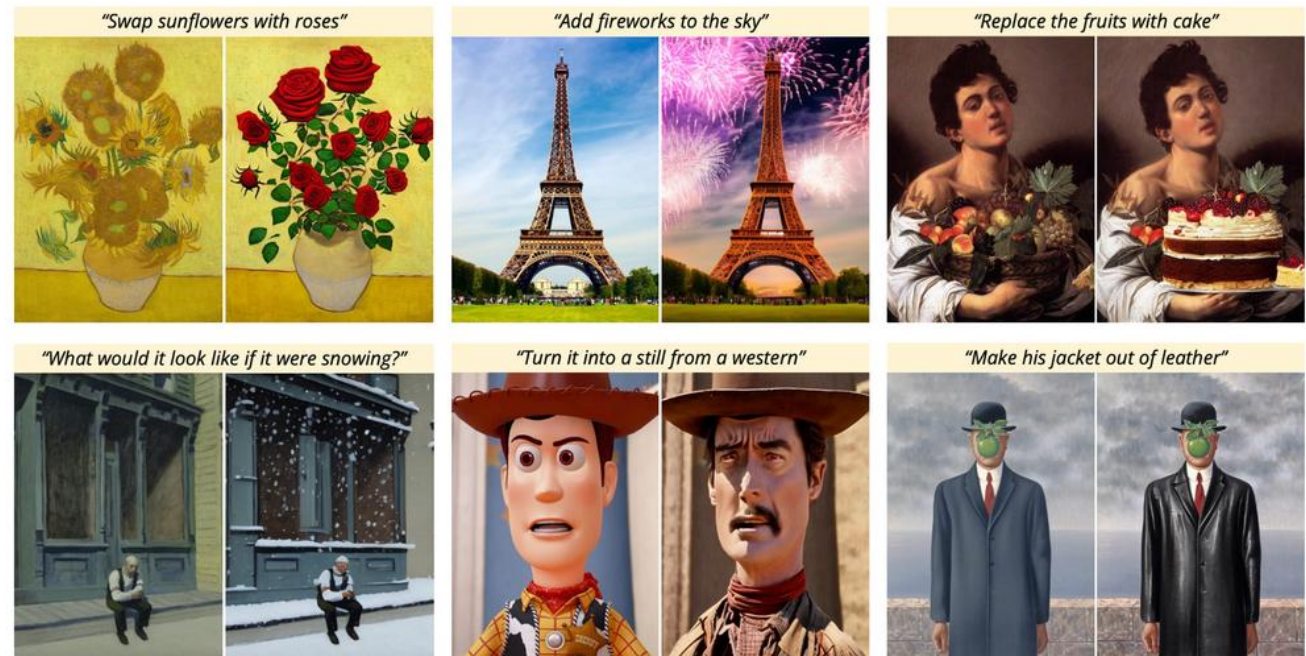
"a fall landscape with a small cottage next to a lake"

Εφαρμογές

34

Image Editing

- Τα μοντέλα διάχυσης μπορούν να χρησιμοποιηθούν για την **επεξεργασία εικόνων**, διατηρώντας το φωτορεαλισμό και την ευκρίνειά τους.
 - Αλλαγή εκφράσεων προσώπου
 - Προσθήκη/αφαίρεση αντικειμένων
 - Αλλαγή χρωμάτων κ.ά.
- Γνωστά μοντέλα:
 - *DreamBooth*
 - *Stable Diffusion (+ ControlNet)*
 - *Instruct Pix2Pix*



Εφαρμογές

35

Image Super-Resolution

- Στο task του Image Super-Resolution θέλουμε να αυξήσουμε την ανάλυση μιας εικόνας χαμηλής ανάλυσης.
- Ιδιαίτερα χρήσιμο για τη βελτίωση της ανάλυσης παλιών φωτογραφιών και φωτογραφιών από πλάνα ασφαλείας.
- Οι Ho et al. 2022 πρότειναν μια αρχιτεκτονική πολλαπλών σταδίων, στην οποία μια σειρά από μοντέλα διάχυσης χρησιμοποιούνται για να αυξήσουν διαδοχικά την ανάλυση των εικόνων.
- Ξεκινώντας με ένα βασικό μοντέλο, το οποίο παράγει μια εικόνα σε χαμηλή ανάλυση, ακολουθούν upsampling μοντέλα, τα οποία αυξάνουν προοδευτικά την ανάλυση και το επίπεδο των λεπτομερειών.

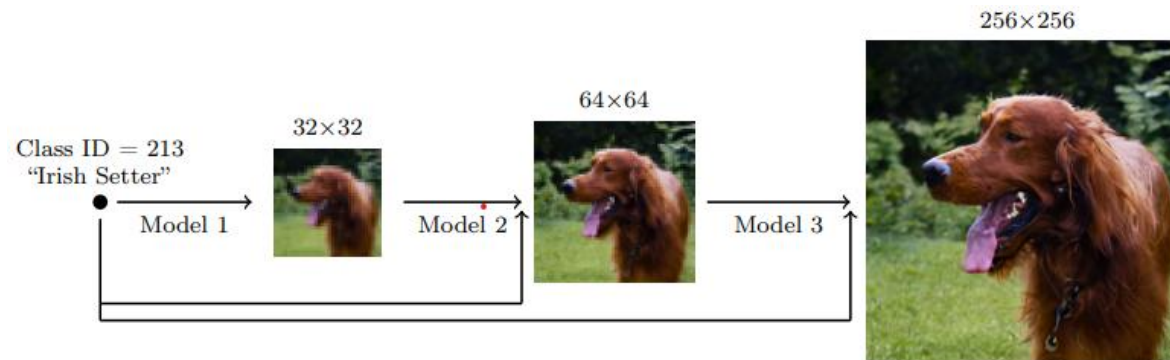


Image-to-Image Translation

- Σκοπός του Image-to-Image Translation είναι η μεταφορά εικόνων από ένα αρχικό πεδίο (source domain) σε ένα άλλο (target domain), διατηρώντας την αναπαράσταση του αρχικού περιεχομένου.
- Πληθώρα εφαρμογών
- **Depth-to-Image Generation**
 - Μετατροπή χάρτη βάθους (depth map) σε εικόνα.
 - Χάρτης βάθους → Εικόνα σε γκριζα κλίμακα. Τα μαύρα αναπαριστούν περιοχές μεγάλου και τα άσπρα περιοχές χαμηλού βάθους.

Source Domain



Target Domain



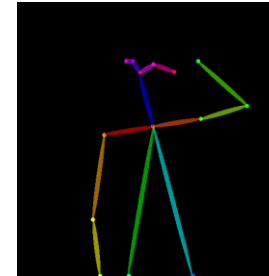
Image-to-Image Translation

- **Pose-to-Image Generation**
 - Μετατροπή εικόνας σκελετού (skeleton image) σε εικόνα.

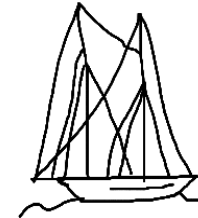
- **Sketch-to-Image Generation**
 - Μετατροπή εικόνας σκίτσου (sketch image) σε εικόνα.

- **Semantic Segmentation-to-Image Generation**
 - Παραγωγή εικόνων από χάρτες σημασιολογικής κατάτμησης (segmentation maps).

Source Domain



Target Domain

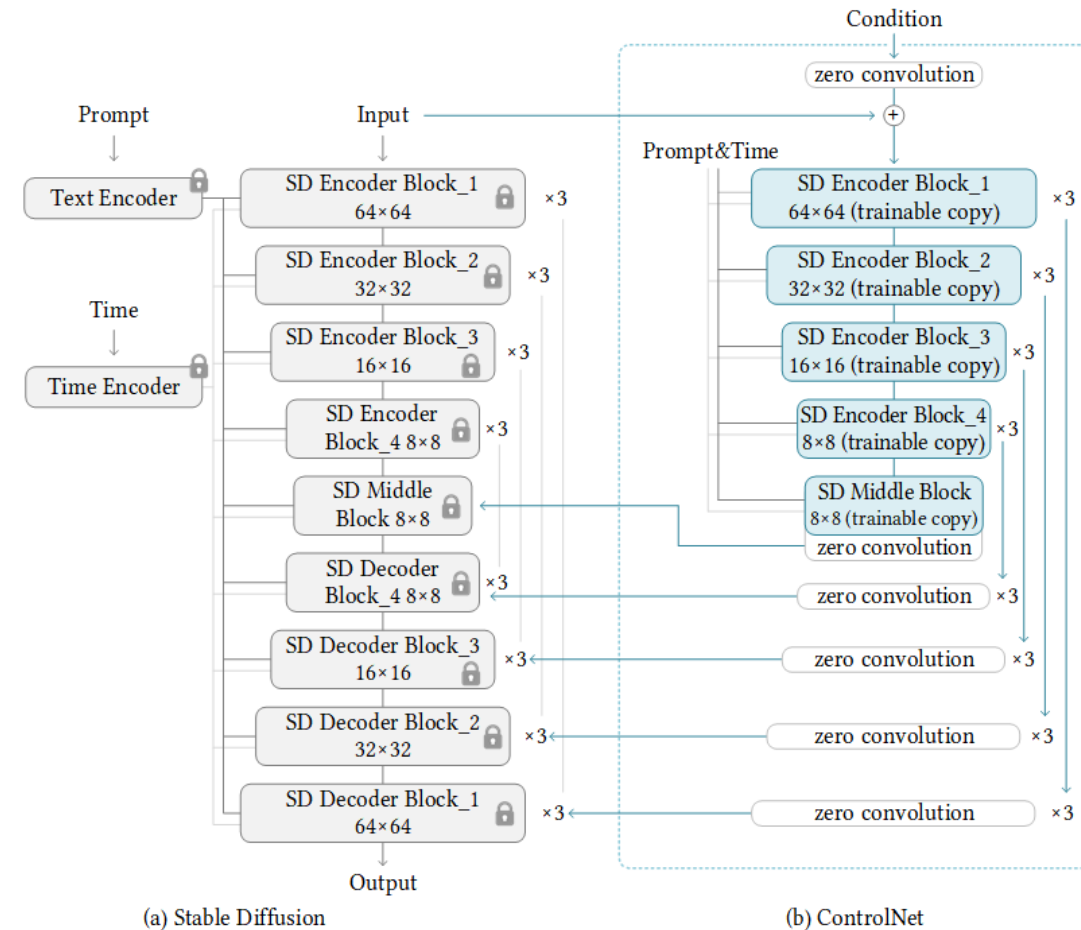


Εφαρμογές

38

Image-to-Image Translation - ControlNet

- Adapter-like δίκτυο, το οποίο λειτουργεί παράλληλα με τη βασική αρχιτεκτονική του Stable Diffusion.
- Επιτρέπει την καθοδήγηση της σύνθεσης βάσει εξωτερικών συνθηκών.
- Κατά την εκπαίδευση το Stable Diffusion παραμένει frozen και ανανεώνονται μόνο τα βάρη του ControlNet.



3D Content Generation

Text-to-3D Shape Generation

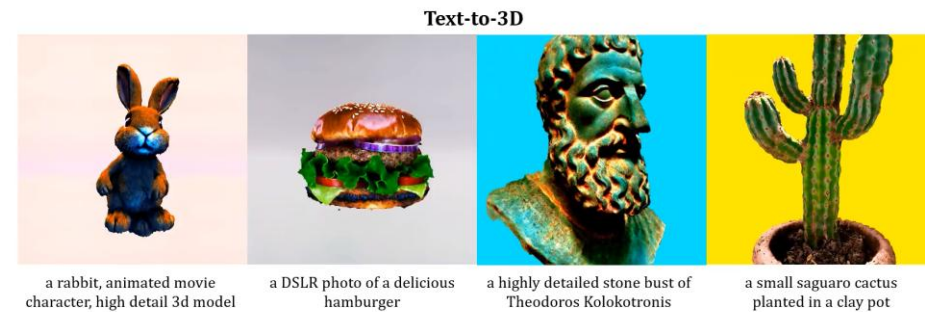


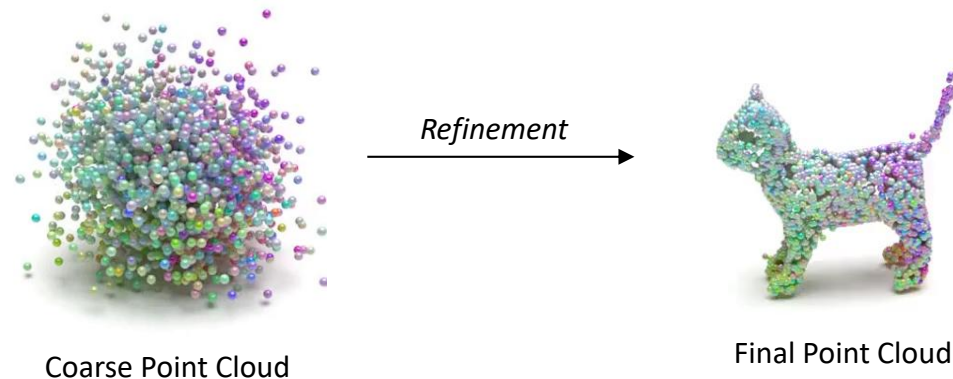
Image-to-3D Scene Generation



3D Content Generation

□ Text-to-Point Cloud Generation

- **Point Cloud:** σύνολο από σημεία στον 3D χώρο, όπου κάθε σημείο έχει συντεταγμένες x, y, z και αναπαριστά ένα συγκεκριμένο σημείο στην επιφάνεια ενός αντικειμένου ή ενός περιβάλλοντος.
- Η αρχιτεκτονική LION της NVIDIA χρησιμοποιεί diffusion models για να παράξει το αρχικό point cloud από τυχαίο θόρυβο και έπειτα κάνει refinement για την παραγωγή του τελικού point cloud.



Text or Image-to-Video Generation

- Τα μοντέλα διάχυσης μπορούν να χρησιμοποιηθούν για την παραγωγή βίντεο από κειμενικές περιγραφές ή εικόνες.
- Γνωστά μοντέλα:
 - *Stable Video Diffusion*
 - *Sora*
 - *ModelScopeT2V*
 - *CogVideoX*

Source Image



Generated Video



Παραγωγή με χρήση του Stable Video Diffusion.

Text Prompt

*“Confident teddy bear surfer
rides the wave in the
tropics”*

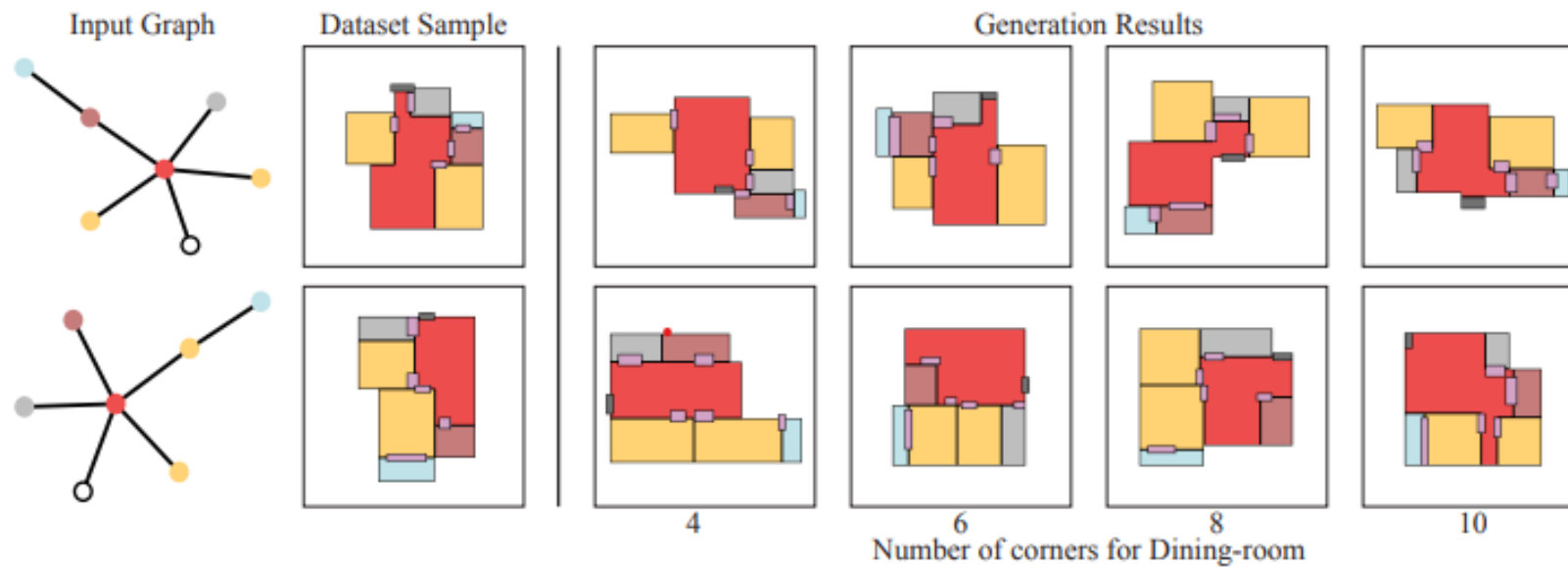
Generated Video



Παραγωγή με χρήση του ModelScopeT2V.

Floorplan Synthesis

- Δημιουργία κατόψεων εσωτερικού χώρου.
- Η διασύνδεση των χώρων αναπαρίσταται μέσω γράφων.
- Ικανότητα ορισμού γωνιών ανά χώρο.



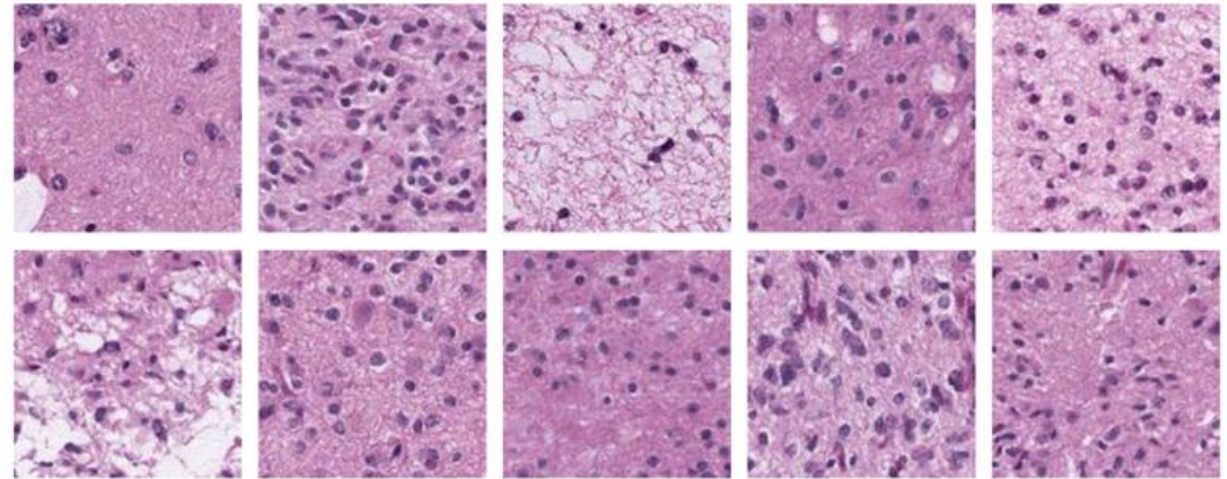
Biomedical Domain

- ▣ Τα παραγωγικά μοντέλα έχουν παρουσιάζουν σημαντική επίδραση στον τομέα της ιατρικής απεικόνισης (medical imaging), παρέχοντας εργαλεία που διευκολύνουν το έργο των ιατρών και των ασθενών.
- ▣ Η διαδικασία συλλογής απεικονιστικών δεδομένων είναι συχνά πολύπλοκη και χρονοβόρα.
 - Εξειδικευμένα πρωτόκολλα και περιορισμοί.
 - Έλλειψη εξειδικευμένων επαγγελματιών για την επισήμανση (annotation) των δεδομένων.
 - Ζητήματα ιδιωτικότητας και ανάγκη για ρητή συναίνεση των ασθενών για διάθεση των προσωπικών τους δεδομένων.
 - Έντονη ανισορροπία στις κλάσεις των δεδομένων, λόγω της φύσης ορισμένων παθολογιών.
- ▣ Τα παραγωγικά μοντέλα προσφέρουν λύσεις στα προβλήματα αυτά καθώς:
 - Δημιουργούν ρεαλιστικά συνθετικά δεδομένα ιατρικών απεικονίσεων.
 - Αίρουν το ζήτημα της προστασίας της ιδιωτικότητας, μέσω δεδομένων που δεν βασίζονται σε πραγματικούς ασθενείς.

Biomedical Domain

■ Synthetic Data Generation

- Έχει αποδειχθεί (Chen et al., 2021) ότι η συνδυαστική χρήση **συνθετικών και πραγματικών δεδομένων** για την εκπαίδευση ταξινομητών σε εικόνες ιστολογίας, **βελτιώνει την απόδοση**.
- Μελέτη η οποία διεξήχθη για την αξιολόγηση των μορφολογικών χαρακτηριστικών συνθετικών και πραγματικών εικόνων σε παθολόγους διαφόρων επιπέδων εμπειρίας (Moghadam et al., 2023), έδειξε ότι οι παθολόγοι δεν μπορούσαν να διακρίνουν τις συνθετικές εικόνες από τις πραγματικές.

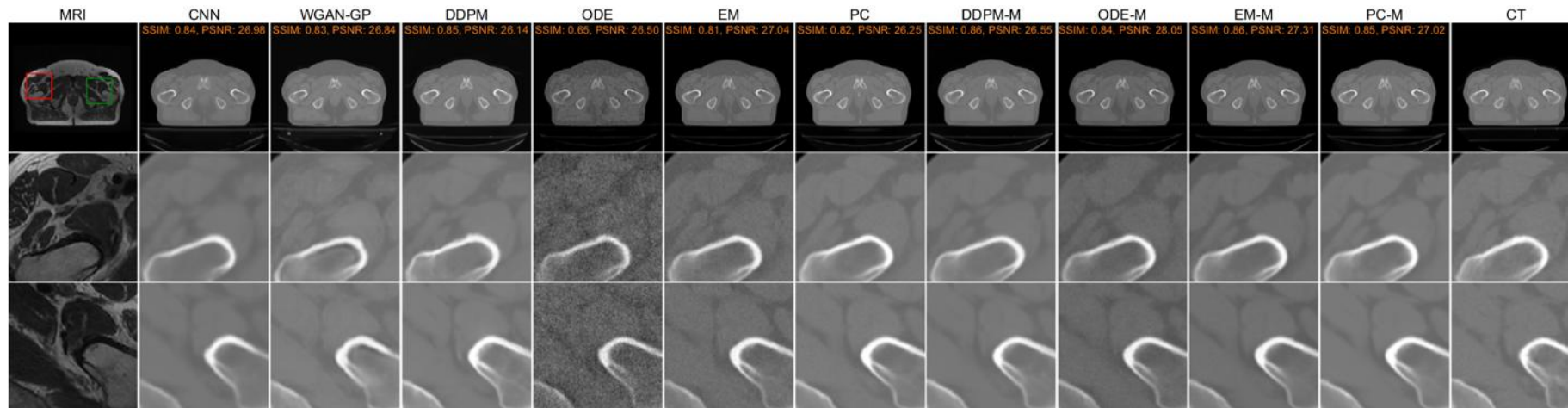


Συνθετικές ιστοπαθολογικές εικόνες που δημιουργήθηκαν από το μοντέλο MFDPM (Moghadam et al., 2023).

Biomedical Domain

Image-to-Image Translation

- Οι πολυτροπικές (multi-modality) ιατρικές εικόνες είναι καθοριστικές για τη διάγνωση και τη θεραπεία.
- Σε αρκετές περιπτώσεις, λείπει κάποια απεικονιστική ακολουθία.
- Τα παραγωγικά μοντέλα έχουν παρουσιάσει ενθαρρυντικά αποτελέσματα στην παραγωγή των απουσιαζουσών απεικονιστικών modalities, π.χ. μετατροπή από *MRI* σε Αξονική Τομογραφία (*CT*).

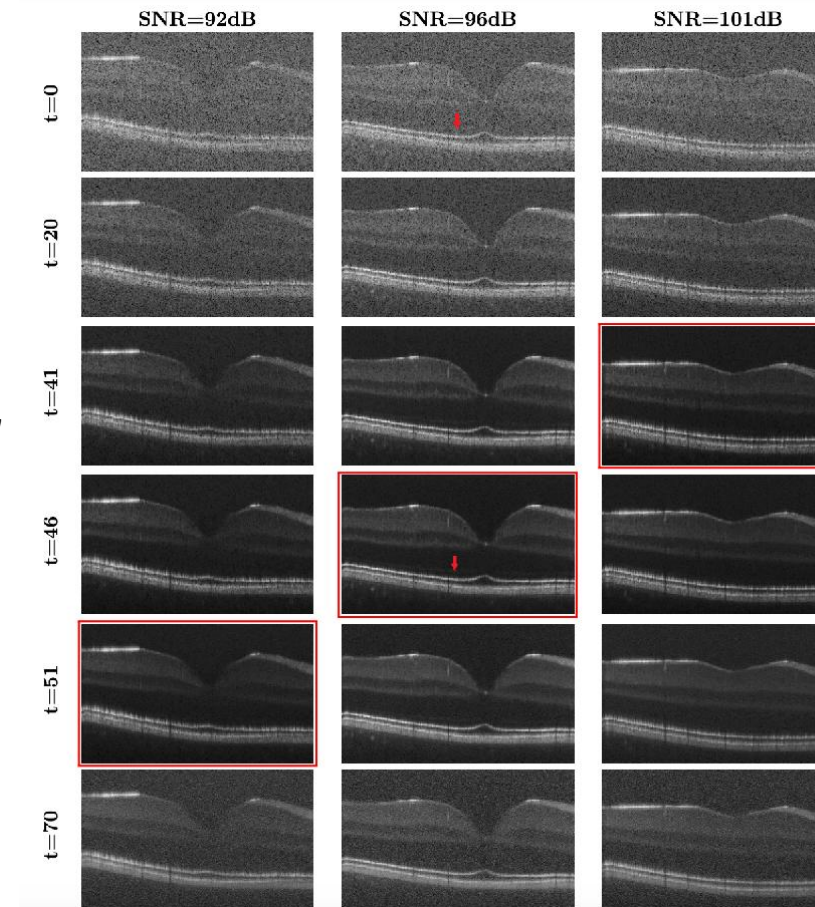


Διαφορετικές μέθοδοι για translation από MRI εικόνες σε CT Scans (Lyu and Wang, 2022).

Biomedical Domain

Image Denoising

- Στην ιατρική απεικόνιση είναι συχνό το φαινόμενο κάποιες εικόνες να επηρεάζονται από θόρυβο κατά τη διάρκεια της απόκτησής τους ή των περαιτέρω σταδίων επεξεργασίας.
- Λόγω τη φύσης τους τα μοντέλα διάχυσης μπορούν να χρησιμοποιηθούν για αποθоруβοποίηση.
- Οι Hu et al. (2022) χρησιμοποίησαν ένα DDPM για την αποθоруβοποίηση εικόνων του αμφιβληστροειδούς χιτώνα (retina) του ματιού.



Αποτελέσματα αποθоруβοποίησης για διάφορα επίπεδα SNR και για διαφορετικές τιμές t . Το καλύτερο για κάθε επίπεδο SNR επισημαίνεται με κόκκινο πλαίσιο (Hu et al., 2022).



Παραγωγή με χρήση του Stable Diffusion