

AI-aided Design?

Processi *text-to-image* per il disegno di architettura

Matteo Flavio Mancini, Sofia Menconero

Abstract

L'Intelligenza Artificiale (AI) sta segnando una svolta in molti campi della vita umana ed è opportuno interrogarsi sulla sua possibilità di utilizzo nei processi di rappresentazione del progetto di architettura. Il contributo presenta una breve digressione sul passato recente delle tecnologie AI al fine di spiegarne il funzionamento, una fotografia sull'attuale stato dell'arte dai processi *text-to-image* a quelli *image-to-3D*, concentrandosi in particolare sulla piattaforma StableDiffusion, oltre a proporre una panoramica sui più recenti studi nel campo del progetto di architettura. La successiva sperimentazione diventa occasione per mostrare le potenzialità dell'AI quanto al processo di co-creazione e alla possibilità di simulare diverse tecniche grafiche, fino alla visualizzazione fotorealistica. D'altro canto, vengono presentati i limiti che, allo stato attuale dello sviluppo, invalidano talvolta i risultati dei processi *text-to-image* per quanto riguarda gli aspetti scientifici della rappresentazione. Le conclusioni propongono una riflessione sulle differenze tra intelligenza umana e artificiale, sul tema dell'autorialità condivisa uomo-macchina e sulle loro conseguenze per il progetto d'architettura.

Parole chiave: intelligenza artificiale, *text-to-image*, disegno di progetto, autorialità, *stablediffusion*.

Introduzione

L'architettura e il disegno di architettura hanno attraversato negli ultimi trent'anni svolte importantissime. Il *first digital turn* [Carpo 2013] ha visto l'introduzione della rappresentazione digitale negli anni Novanta del XX secolo mentre il *second digital turn* [Carpo 2017] si è avviato con la diffusione di algoritmi e *big data* a partire dagli anni '10 del XXI secolo. Dieci anni dopo, stiamo assistendo a un'altra potenziale svolta dovuta a una repentina accelerazione dello sviluppo e della diffusione di strumenti di Intelligenza Artificiale (AI), già in uso nei maggiori studi di architettura come Coop Himmelb(l)au [Prix et al. 2022], Zaha Hadid Architects [Wallish 2022] e Foster + Partners [Tsigkari et al. 2021].

Una branca dell'AI, basata su processi di tipo *text-to-image*, propone soluzioni facilmente accessibili e dedicate alla

creazione di immagini. Si tratta di un modello di *machine learning* che usa un linguaggio naturale descrittivo come input e produce un'immagine basata sull'elaborazione della descrizione fornita. I risultati ottenuti da queste piattaforme sono sorprendenti in termini di corrispondenza ai *prompt* testuali inseriti e di flessibilità delle tecniche grafiche che sono in grado di (ri)produrre.

Partendo dal presupposto che, allo stato attuale, queste AI non hanno alcuna coscienza creativa né un'effettiva capacità di comprendere le regole compositive e proiettive o la spazialità rappresentata nelle immagini, è comunque opportuno interrogarsi sulle loro possibilità di utilizzo nei processi di rappresentazione del progetto di architettura. Con questo obiettivo e tenendo conto delle caratteristiche

intrinseche di questa tecnologia, che verranno espone nei successivi paragrafi, si propone la sperimentazione dell'AI *text-to-image* attraverso la piattaforma open-source *StableDiffusion* per la realizzazione di immagini prospettiche capaci di contribuire alle fasi preliminari di ideazione del progetto.

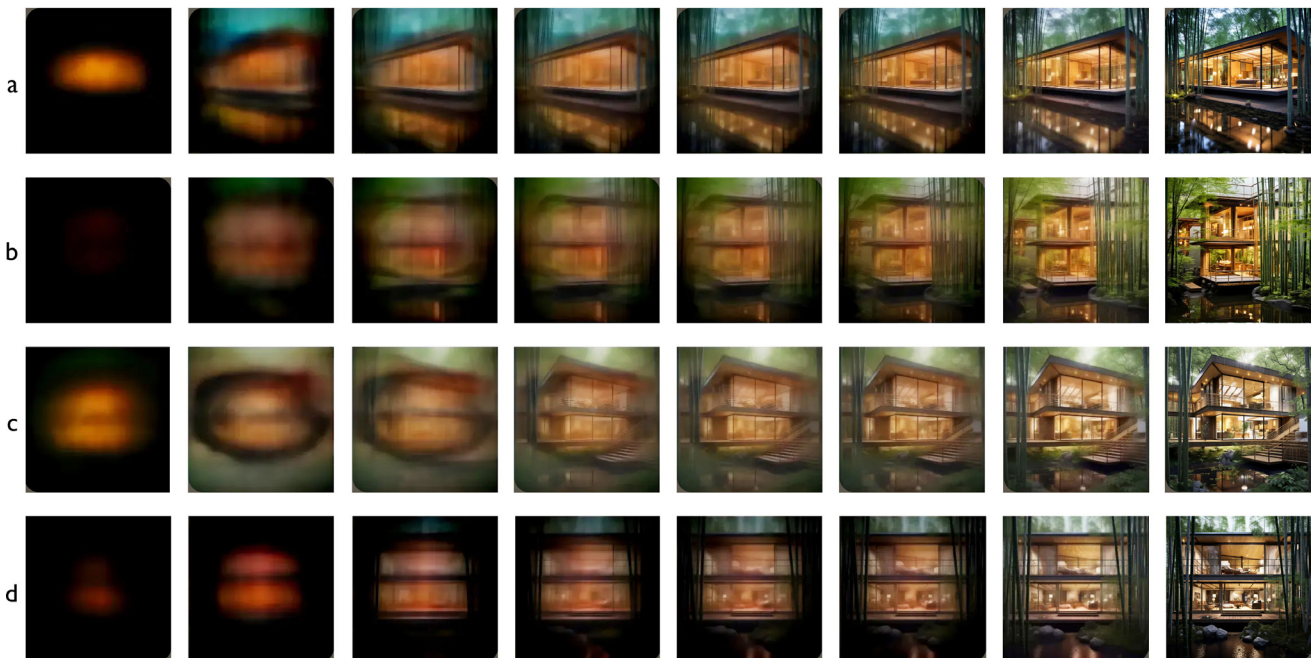
Passato (recente) e presente della generazione di immagini basata sull'AI

I sistemi generativi basati sull'AI che è possibile fruire nel campo dell'architettura e del design sono in rapida evoluzione. Un punto di svolta nella ricerca sulla generazione di immagini è segnato dall'invenzione delle *Generative Adversarial Networks (GAN)* nel 2014 [Goodfellow et al. 2014]. Si tratta di un'architettura di *deep learning* in cui due reti

neurali antagoniste (un generatore e un discriminatore) interagiscono reiterativamente durante l'addestramento al fine di giungere al punto in cui il discriminatore non sia più in grado di distinguere le immagini sintetiche prodotte dal generatore rispetto alle immagini reali immesse come dati di addestramento. Nel 2016 è stata sviluppata un'architettura GAN in grado di generare immagini plausibili da descrizioni testuali dettagliate [Reed et al. 2016] avviando di fatto il sistema AI *text-to-image*. Un ulteriore avanzamento è stato segnato da un più efficiente metodo di apprendimento basato sull'elaborazione del linguaggio naturale, chiamato CLIP (*Contrastive Language-Image Pre-training*) [Radford et al. 2021]. Questo modello di classificazione delle immagini identifica gli oggetti imparando dal testo associato a un'immagine (piuttosto che da etichette assegnate manualmente) ed è stato allenato su 400 milioni di coppie immagini-testo estratte dal web. I modelli CLIP

Fig. 1. Processo di denoising durante la generazione di quattro varianti (a, b, c, d) in Midjourney attraverso il prompt: *a modern Japanese house in a bamboo forest in spring* (elaborazione degli autori).

denoising process



sono in grado di stimare la conformità di un'immagine generata per un *prompt* testuale [Colton et al. 2021]. Un esempio di questo abbinamento è il sistema di generazione di immagini chiamato *VQGAN-CLIP*, che utilizza una rete neurale GAN ancora più potente. I contributi significativi dell'architettura *VQGAN-CLIP* riguardano la qualità visiva sia nella generazione che nella manipolazione delle immagini, la fedeltà semantica tra testo e immagine generata, l'efficienza dovuta al fatto che il metodo non richiede ulteriore addestramento oltre ai modelli pre-allenati e il valore dello sviluppo e della scienza aperta [Crowson et al. 2022, p. 2]. Successivamente, i sistemi basati su GAN sono stati sostituiti con i *diffusion model*, ovvero modelli probabilistici di *machine learning* addestrati a eliminare il rumore delle immagini precedentemente introdotto, imparando a invertire il processo di diffusione [Dhariwal et al. 2021]. L'addestramento di questi modelli li rende in grado di utilizzare i metodi di *denoising* per sintetizzare nuove immagini, prive di rumore, da input casuali (fig. 1).

Alcune applicazioni dell'AI fruibili nel campo dell'architettura e del design sono le seguenti:

- *text-to-image*, la più diffusa operazione di generazione di immagini attraverso una descrizione testuale, spesso associata ad altre funzionalità;
- *image-to-image*, per la trasformazione di un'immagine input in modo che corrisponda alle caratteristiche di un'immagine di destinazione, può essere usata per trasferire uno stile, per modificare o rimuovere oggetti dalle immagini (*inpainting*), per trasformare a colori un'immagine in bianco e nero, per aumentare la risoluzione di un'immagine (*upscaling*);
- *text* o *image-to-video*, per creare video da un *prompt* testuale (ad esempio *Make-a-Video* [Singer et al. 2022], o *CogVideo* [Hong et al. 2022]) oppure per creare un'animazione grazie al montaggio di immagini generate attraverso l'*image-to-image* (ad esempio *Deform*) con effetti simili a un video in *stop-motion*.
- *text* o *image-to-3D*, per generare modelli 3D da un *prompt* testuale (ad esempio *Point-E* per generare nuvole di punti [Nichol et al. 2022], *Shape-E* per mesh texturizzate [Jun et al. 2023]), oppure la generazione dei modelli 3D può avvenire a partire da un'immagine (ad esempio *Kaedim*).

L'incredibile recente diffusione dell'AI *text-to-image* deriva dall'attivazione di alcune piattaforme con interfacce semplici da utilizzare, anche da parte di fruitori non esperti, come *DALL-E 2*, *Midjourney* e *StableDiffusion*.

Le principali piattaforme per l'AI *text-to-image*

DALL-E è la prima tra le tre piattaforme a essere stata presentata nel gennaio 2021 (l'attuale versione 2 è di aprile 2022) da *OpenAI* [Ramesh et al. 2021], gli stessi sviluppatori di *ChatGPT*. La piattaforma, fruibile online in abbonamento, propone quattro funzioni: la generazione di immagini realistiche e artistiche da una descrizione testuale che può combinare concetti, attributi e stili (fig. 2); l'*outpainting*, ovvero l'espansione dell'immagine oltre i margini originali attraverso la creazione di una nuova composizione; l'*inpainting*, tramite il quale è possibile modificare porzioni di immagine aggiungendo o eliminando degli oggetti attraverso una descrizione testuale e mantenendo coerente il resto della scena; la generazione di variazioni ispirate a un'immagine di input.

Midjourney, rilasciata il 12 luglio 2022, è attualmente giunta alla versione 5.2 con notevoli miglioramenti rispetto agli esordi in termini di aderenza al *prompt* e fotorealismo (fig. 3) e conta, dopo un anno, oltre 15 milioni di utenti [1].

Fig. 2. Immagine generata con DALL-E 2 attraverso il *prompt*: *a modern building on a crowded street at sunset* (elaborazione degli autori).



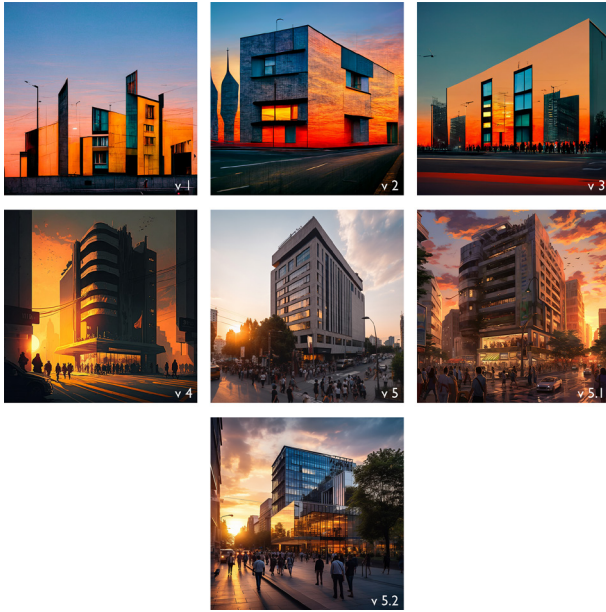


Fig. 3. Confronto tra diverse immagini generate a partire dallo stesso testo (prompt: a modern building on a crowded street at sunset) relative a diverse versioni di Midjourney (elaborazione degli autori).

Come *DALL-E*, è fruibile online in abbonamento. Le tre principali attività generative su *Midjourney* sono: un'immagine a partire da un *prompt* testuale; una descrizione a partire da un'immagine; un'immagine di sintesi a partire da due fino a cinque immagini in input. *Midjourney* (come *StableDiffusion*) permette di utilizzare anche un *prompt* negativo nel caso non si vogliono specifici elementi nell'immagine generata.

StableDiffusion, rilasciata ad agosto 2022, è l'unica delle tre piattaforme a essere open-source e si basa su un *diffusion model* chiamato *latent diffusion model* [Rombach et al. 2022]. L'attuale versione beta XL è disponibile solo online in abbonamento, mentre le precedenti versioni possono anche essere installate in locale gratuitamente. *StableDiffusion* supporta la generazione di immagini attraverso l'uso di un *prompt* di testo che descrive gli elementi da includere o escludere dall'output (fig. 4), l'*inpainting* e l'*outpainting*, la generazione *image-to-image* e l'*upscaling*. È inoltre possibile associare a *StableDiffusion* delle estensioni come *ControlNet*, che genera variazioni di un'immagine

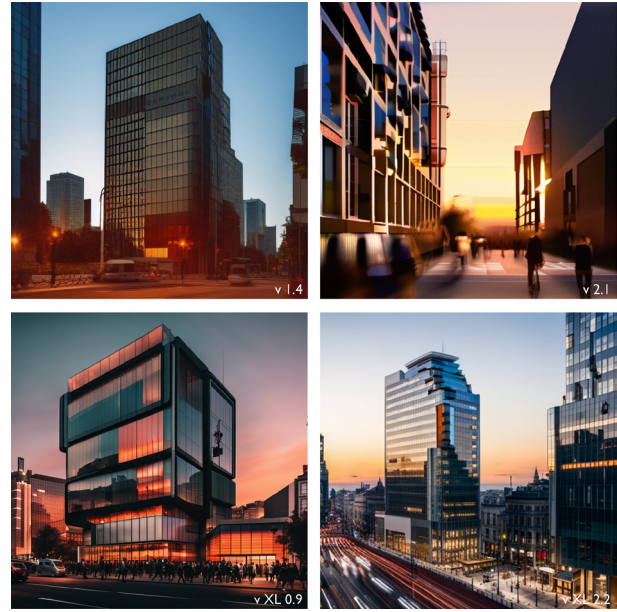


Fig. 4. Confronto tra diverse immagini generate a partire dallo stesso testo (prompt: a modern building on a crowded street at sunset) relative a diverse versioni di *StableDiffusion* (elaborazione degli autori).

di input attraverso descrizioni testuali, e *Deform*, che attraverso la funzione *image-to-image* genera una serie di immagini, applicando piccole trasformazioni, e le cuce insieme per creare un video.

Il vantaggio più grande di *StableDiffusion* rispetto alle altre piattaforme è la possibilità che gli utenti finali possano implementare un addestramento aggiuntivo (*fine-tuning*) per ottimizzare gli output di generazione in modo che corrispondano a casi d'uso più specifici. Ad esempio, negli studi di architettura dove l'AI è entrata a far parte del processo creativo, la rete neurale viene allenata con immagini mirate del repertorio progettuale dello studio al fine di ottenere risultati più in linea con il linguaggio architettonico e grafico.

Dunque, a differenza delle precedenti piattaforme, *StableDiffusion* permette una maggiore libertà di utilizzo in termini di personalizzazione del processo generativo, per questo motivo è stata scelta per la successiva sperimentazione, associata all'estensione *ControlNet* [Zhang et al. 2023] la quale migliora il controllo degli output.

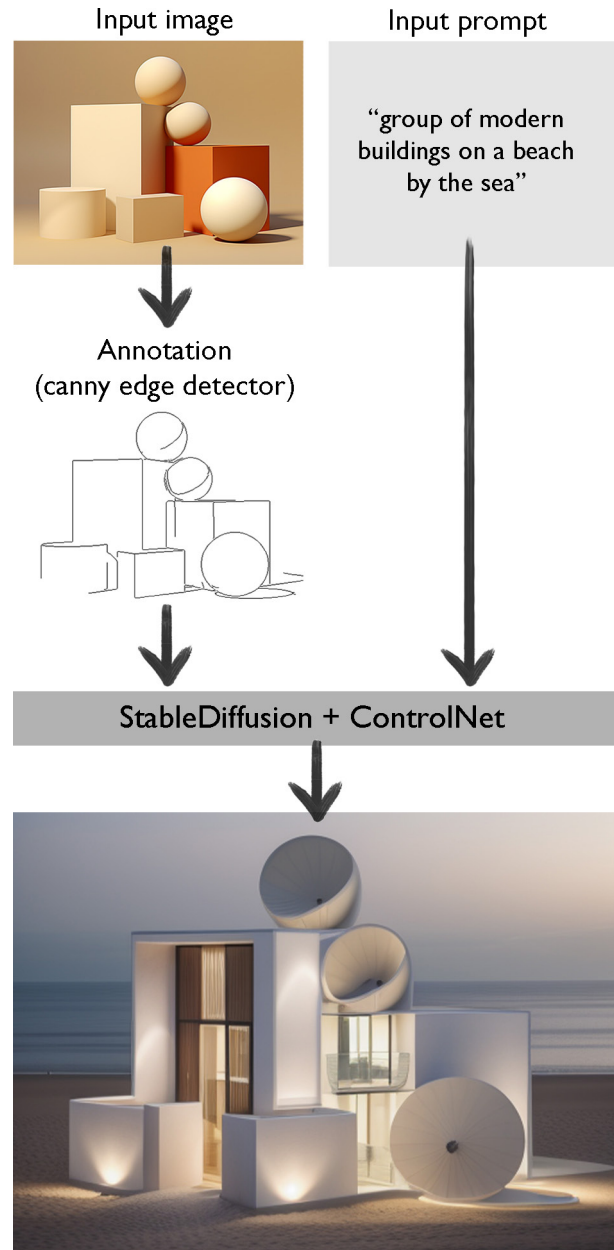
Quest'ultima è una struttura di rete neurale progettata per gestire modelli di diffusione incorporando condizioni aggiuntive: manipolando le condizioni di input dei blocchi riesce a controllare ulteriormente il comportamento generale di un'intera rete neurale. *ControlNet* agisce a partire da un'immagine in input e una descrizione testuale, e permette di ottenere delle immagini che sono variazioni conformi all'input dal punto di vista compositivo ma che seguono anche la descrizione impostata. Il processo a cui i dati sono sottoposti prevede, innanzitutto, la generazione di una mappa basata sull'immagine input (chiamata fase di annotazione o pre-processamento) la quale viene usata dalla rete per generare le varianti con le caratteristiche descritte testualmente (fig. 5).

Avanti e indietro nello spazio latente: tra allenamento e generazione nel modello di diffusione di *StableDiffusion*

Per approcciare correttamente le sperimentazioni che seguiranno nei prossimi paragrafi è importante cercare di comprendere non tanto gli aspetti prettamente tecnico-informatici quanto i processi attuati da questo tipo di AI, nello specifico *StableDiffusion*, nei due momenti distinti dell'allenamento e della generazione, poiché da essi dipendono sia l'uso appropriato che l'interpretazione critica di questa tecnologia.

I modelli di diffusione mutuano dalla termodinamica il concetto di diffusione, cioè il fenomeno per cui le particelle di un fluido si muovono randomicamente all'interno di un altro fluido con diversa concentrazione, fino a raggiungere una nuova condizione di equilibrio. Allo stesso modo, le immagini delle AI durante la generazione sembrano progressivamente emergere dal caos del rumore digitale. Il principio della diffusione viene usato sia in fase di allenamento (*forward diffusion*), che in fase di generazione (*reverse diffusion*). In *StableDiffusion*, entrambi questi processi avvengono nel *latent space*, uno spazio numerico/informativo in cui le immagini vengono tradotte in tensori (matrici a più dimensioni) per lavorare su una loro versione compressa, più leggera del *pixel space* iniziale delle immagini. Anche i testi che descrivono le immagini subiscono una simile traduzione e compressione. L'analogia tra rappresentazione latente dei testi e delle immagini è importante

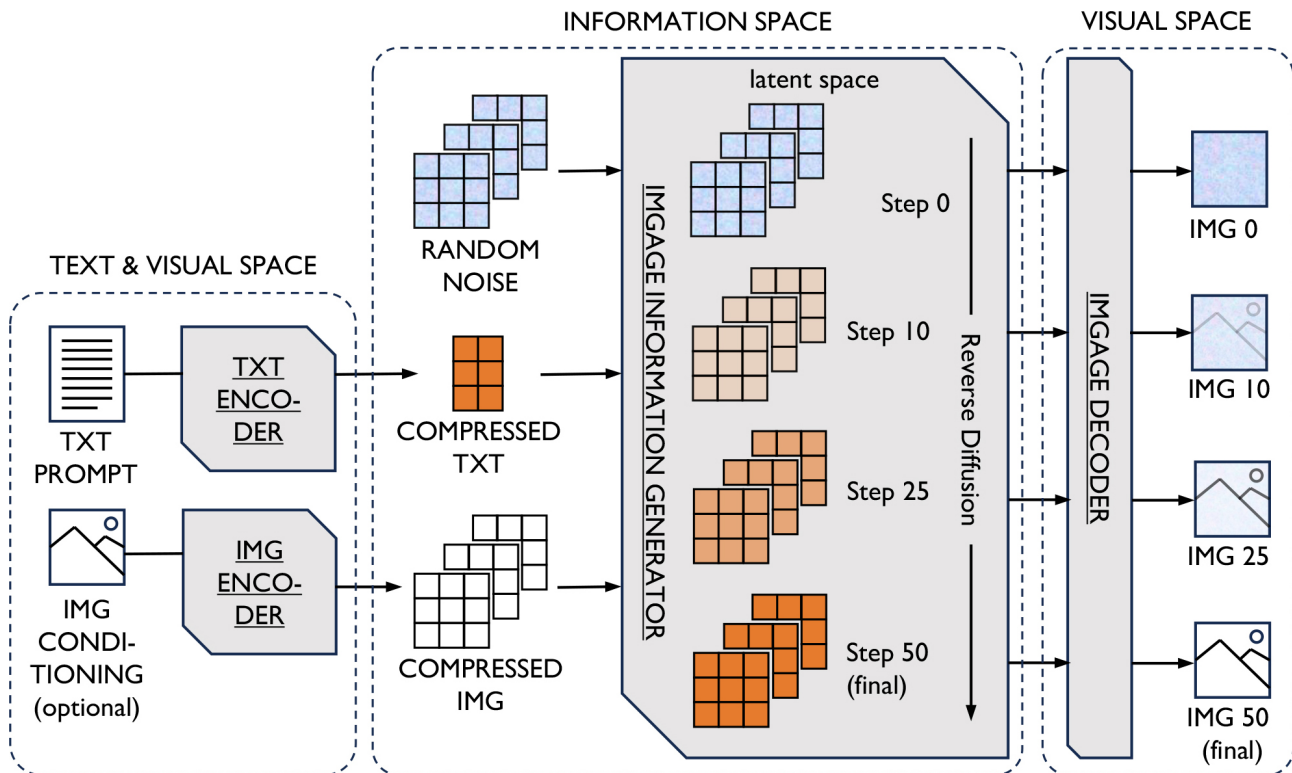
Fig. 5. Schema della generazione di un'immagine attraverso *StableDiffusion* con il condizionamento (*canny edge*) di *ControlNet* (elaborazione degli autori).



perché fa capire come le AI non immagazzinino e non elaborino raccolte di sillabe, parole o porzioni di immagini ma operino su rappresentazioni numeriche astratte delle caratteristiche delle immagini, degli oggetti rappresentati, delle possibili situazioni ambientali e delle varie tecniche e stili. Il *latent space* può essere immaginato come il luogo dove l'AI memorizza, in fase di allenamento, e da cui richiama, in fase di generazione, le proprie 'conoscenze'. L'allenamento di queste AI non è progressivo nel tempo ma avviene prima della pubblicazione, pertanto le loro "conoscenze" sono statiche e aggiornate periodicamente con il progredire delle versioni. Ad esempio, *StableDiffusion* è allenato sul dataset open-access LAION 5B, composto da 5 miliardi di coppie immagine-didascalia il cui contenuto è esplorabile, a partire dall'immissione di un *prompt*

testuale, attraverso un apposito portale [2]. La consultazione del dataset di allenamento permette di farsi un'idea sulle corrispondenze tra termini e immagini e, quindi, su cosa sia possibile aspettarsi dai risultati delle elaborazioni: una ricerca che non restituisca risultati coerenti indica che non potranno essere generate delle immagini che corrispondano alle aspettative per quell'input testuale. In fase di addestramento, le immagini del database vengono elaborate introducendo pattern casuali di rumore di diversa intensità. Le immagini così elaborate, insieme alle didascalie corrispondenti, vengono sottoposte all'AI per allenarla a individuare il tipo di pattern adottato, la quantità di rumore introdotta e a rimuoverne entrambi per migliorare la qualità delle immagini. In questo modo, attraverso il processo del *forward diffusion*, l'AI apprende contemporaneamente sia

Fig. 6. Schema delle fasi del processo di generazione di un'immagine attraverso il modello di diffusione adottato da StableDiffusion (elaborazione degli autori).



come ottenere immagini prive di rumore sia le corrispondenze tra immagini e testi.

Quanto appreso viene utilizzato da *StableDiffusion* per sviluppare un processo generativo che parte e termina in uno spazio in cui i dati (testi e immagini) sono adatti alla percezione umana, attraversa uno spazio puramente informativo (spazio latente) in cui i dati sono rappresentati da *token* (testi) e *tensori* (immagini).

Il processo generativo può essere suddiviso in tre blocchi fondamentali (fig. 6). Il primo prevede la compressione e trasformazione in informazioni numeriche attraverso *encoder* (reti neurali appositamente allenate) dei dati input inseriti per condizionare la generazione dell'immagine. In *StableDiffusion*, grazie all'estensione *ControlNet*, gli input testuali possono essere integrati da condizionamenti grafici opzionali. Nel secondo blocco, attraverso il processo di *reverse diffusion*, avviene l'elaborazione degli input in relazione alle conoscenze note. Tale processo è reiterativo e passa attraverso diversi step di *denoising* per affinare la corrispondenza tra gli input immessi e l'immagine generata. In questa fase l'elaborazione avviene a livello di informazioni numeriche e non c'è alcuna elaborazione grafica di immagini. Quest'ultima avviene nel terzo blocco, dove le rappresentazioni numeriche vengono tradotte da una rete neurale con funzione di *decoder* in immagini visualmente percepibili [Rombach et al. 2022].

Studi relativi all'AI applicata al progetto di architettura

Alcuni studi che riguardano l'AI nell'ambito della progettazione architettonica sono incentrati nell'evidenziare le potenzialità e i limiti della tecnologia. Nella maggior parte dei casi il potenziale è riscontrato come supporto nel processo creativo [Jaruga-Rozdolska 2022; Paananen et al. 2023]. Tra le altre potenzialità sono citate l'abilità di poter immaginare forme astratte, re-immaginare l'architettura biomimetica, rivisitare l'architettura tradizionale e visualizzare avanzamenti fotorealistici a partire da schizzi architettonici. I limiti individuati sono relativi alla possibilità di controllo e personalizzazione dei processi, alla scarsa considerazione per quanto riguarda gli aspetti di fattibilità strutturale, all'eventuale incoerenza stilistico-architettonica dei risultati generati [Hegazy et al. 2023]. I casi studio relativi al progetto di architettura su cui l'AI è stata applicata riguardano la fase ideativa, la generazione di schizzi con specifici stili grafici, l'aggiunta di persone e oggetti in immagini esistenti, la combinazione di

varie parti di immagini in una composizione coerente, la variazione di un'immagine iniziale, la variazione dello stile grafico di un'immagine esistente, il disegno planimetrico, il design di esterni e interni, la creazione di texture, il progetto urbano [Ploenning et al. 2022; Yildirim 2022].

Uno studio didattico riguarda l'integrazione delle tecniche di AI alle tecniche tradizionali in un corso di rappresentazione del design al primo anno universitario, dove gli autori hanno notato un miglioramento delle capacità interpretative e compositive degli studenti [Tong et al. 2023]. Agli studenti era stato chiesto di creare una composizione di solidi e di disegnare a mano proiezioni ortogonali e assonometria isometrica; poi di generare una serie di immagini con *Midjourney* attraverso alcune parole chiave; infine di combinare le due produzioni precedenti mediante varie tecniche.

Potenzialità dell'AI text-to-image per il disegno di architettura

Per sperimentare il possibile contributo dell'AI nella fase preliminare del progetto, il momento in cui la rappresentazione contribuisce all'ideazione e alla prefigurazione, si è deciso di lavorare sia riguardo alla definizione dell'idea che alla sua visualizzazione.

Sono stati ipotizzati tre diversi input grafici: due viste prospettiche esterne di un modello tridimensionale volumetrico e uno schizzo al tratto di un interno, tutti volutamente privi di caratterizzazioni se non quelle minime indispensabili per la definizione spaziale e l'impostazione dell'inquadratura. Questi input grafici, grazie all'estensione *ControlNet*, hanno il compito di inserire nel processo generativo l'impostazione morfologica generale del progetto mentre gli input testuali vengo utilizzati per descrivere le tecniche grafiche desiderate ed eventuali caratteristiche delle architetture in fatto di materiali, contesto e ulteriori caratteristiche stilistiche che si desidera inserire. I risultati di queste prime sperimentazioni dimostrano la notevole flessibilità dell'AI nel (ri)creare tecniche grafiche diverse, che variano dal disegno a lapis, alle matite colorate, all'acquerello, con una notevole capacità di integrazione di elementi di contesto sia naturali che artificiali. Contemporaneamente, l'aggiunta da parte dell'AI di elementi di dettaglio quali trame, bucature e materiali, contribuisce all'avanzamento dell'ideazione in un processo in cui si può ipotizzare che alcuni di questi elementi possano essere effettivamente inseriti nel prosieguo del progetto, in uno scambio uomo-macchina reiterato (fig. 7).

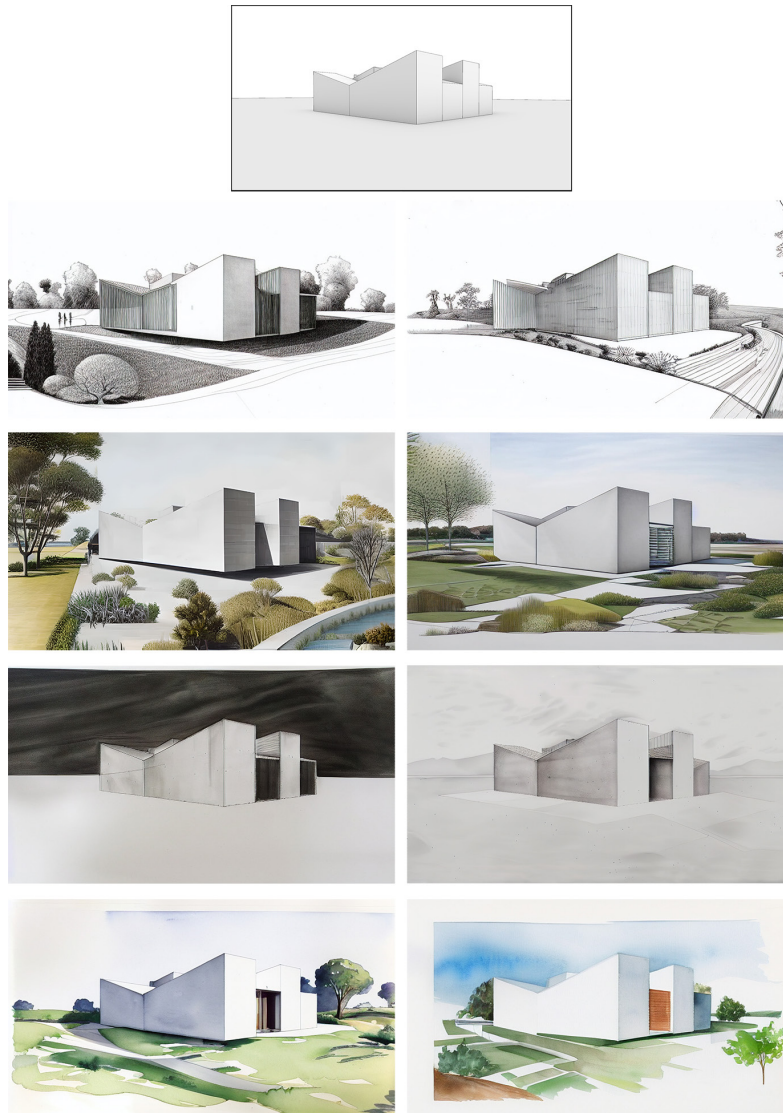


Fig. 7. Immagini generate con StableDiffusion per simulare diverse tecniche grafiche. Dall'alto verso il basso: lapis, matite colorate, acquerello monocromatico e a colori. Prompt: linear, exterior view, contemporary architecture, highly detailed architecture, large windows, concrete, architectural drawings, technical drawings, [tecnica grafica desiderata], line drawings, working drawings, architectural sketches, conceptual style, abstract (elaborazione degli autori).

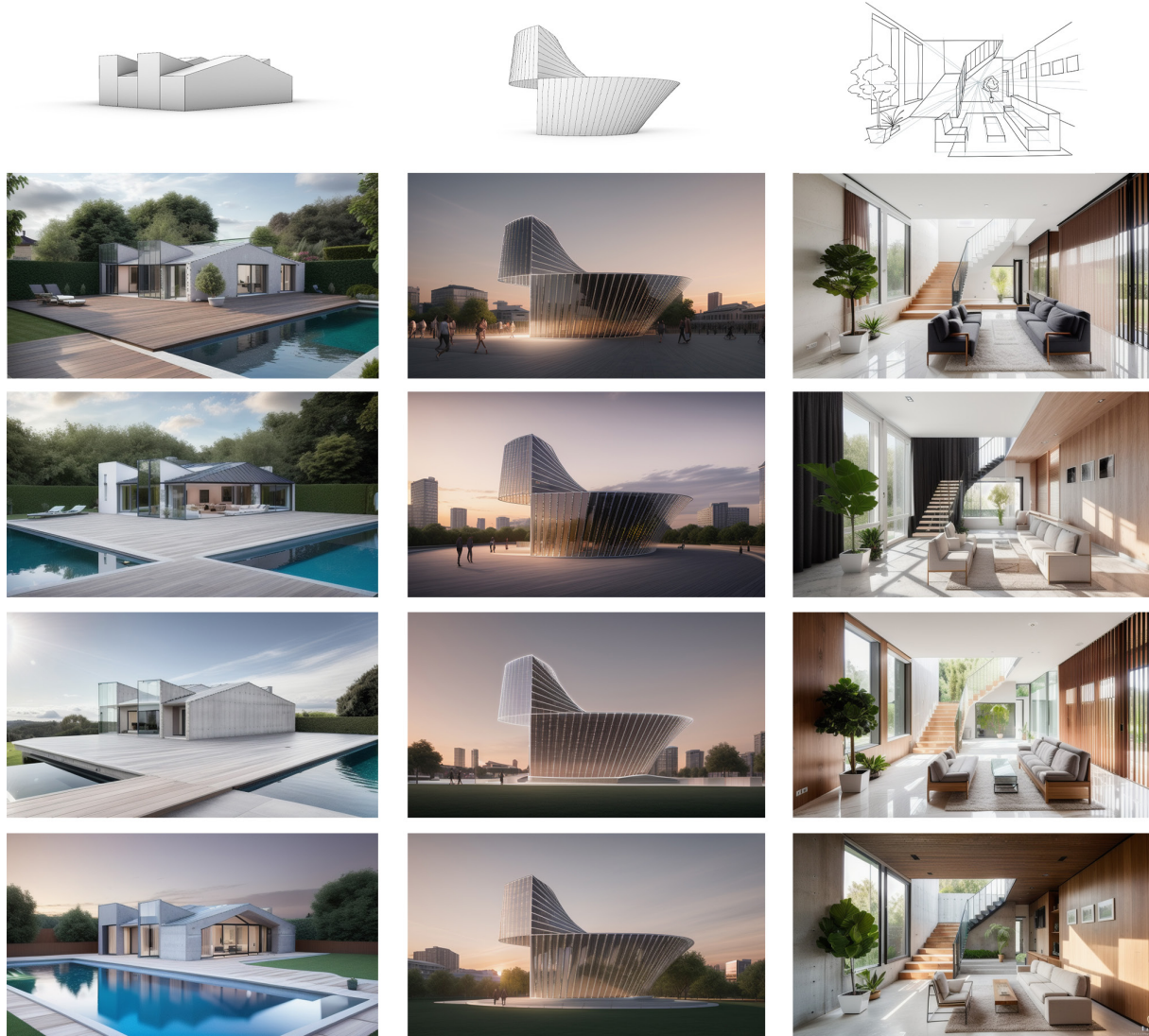


Fig. 8. Immagini generate con StableDiffusion per visualizzazioni fotorealistiche. A sinistra, viste esterne, prompt: exterior home view, concrete walls and roof, large glass windows, small rectangular swimming pool, garden, garden furniture, clouds. Al centro, viste esterne con superfici curve, prompt: a pavilion in a contemporary architecture style, covered with reflective panels, surrounded by a round pool with people and trees. A destra, viste interne, prompt: home interior view, modern architecture, large glass windows with curtains, timber framing, wood flooring, concrete ceiling, steel staircase, large sofa with pillows, armchair, coffee table with flowerpot, carpet, plants, lamp, minimalist style furniture, sunlight from windows, daylight (elaborazione degli autori).

Midjourney

StableDiffusion +
RealisticVision v4

Il possibile contributo in termini di definizione dell'idea attraverso la rapida generazione di varianti è più evidente se si richiede all'AI di produrre immagini fotorealistiche. In questo caso si apprezza maggiormente la capacità di proporre variazioni a partire da quanto richiesto tramite il *prompt* testuale. Le sperimentazioni condotte sulle viste esterne mostrano la varietà di materiali e interpretazioni dei semplici schemi volumetrici proposti come input, nonché l'abilità nella creazione di contesti di ambientazione (fig. 8). Analogamente, le sperimentazioni basate su uno schizzo digitale al tratto di un ambiente interno mettono in luce la capacità di accostamento cromatico e dei materiali ma anche la propensione ad aggiungere elementi quali tende e sopralzi del pavimento. Compaiono anche elementi di piccole dimensioni, quali punti luce e complementi di arredo. La distribuzione di queste integrazioni appare in linea di massima coerente con l'impostazione generale.

Limiti dell'AI *text-to-image* nella rappresentazione

I limiti indagati nel presente paragrafo [3] riguardano in particolare l'aspetto della rappresentazione architettonica (prospettiva, riflessioni, illuminazione/ombre).

Allo stato attuale, l'AI non ha alcuna coscienza delle regole proiettive sottese a una corretta costruzione prospettica. Se da un lato questa asserzione era deducibile sulla base dei principi teorici dietro la tecnologia, essa trova anche conferma su base sperimentale. Andando ad aggiungere, nel *prompt* testuale, una parte descrittiva riguardante il metodo di rappresentazione (*central perspective*) [4] si giunge a dei risultati in cui la prospettiva centrale è presente solo in una parte delle immagini generate (fig. 9). Andando successivamente ad analizzare l'impianto prospettico di due delle precedenti immagini generate, si osserva che le linee di fuga delle piastrelle quadrate del pavimento (rette orizzontali perpendicolari al piano di quadro) non individuano un punto univoco di convergenza (fig. 10). Inoltre, tracciando le diagonali delle piastrelle quadrate dai due estremi visibili nelle immagini, si nota che le intersezioni intermedie non sono perfettamente allineate alle diagonali. Dunque, le prospettive sono perettivamente efficaci ma non sono proiettivamente corrette. I risultati della sperimentazione prospettica fanno supporre che l'AI non sia stata addestrata a riconoscere correttamente i diversi metodi della rappresentazione.

Fig. 9. Immagini generate attraverso il *prompt*: *central perspective, home interior view, floor with regular dark square tiles, modern architecture, minimalist style furniture, daylight* (elaborazione degli autori).

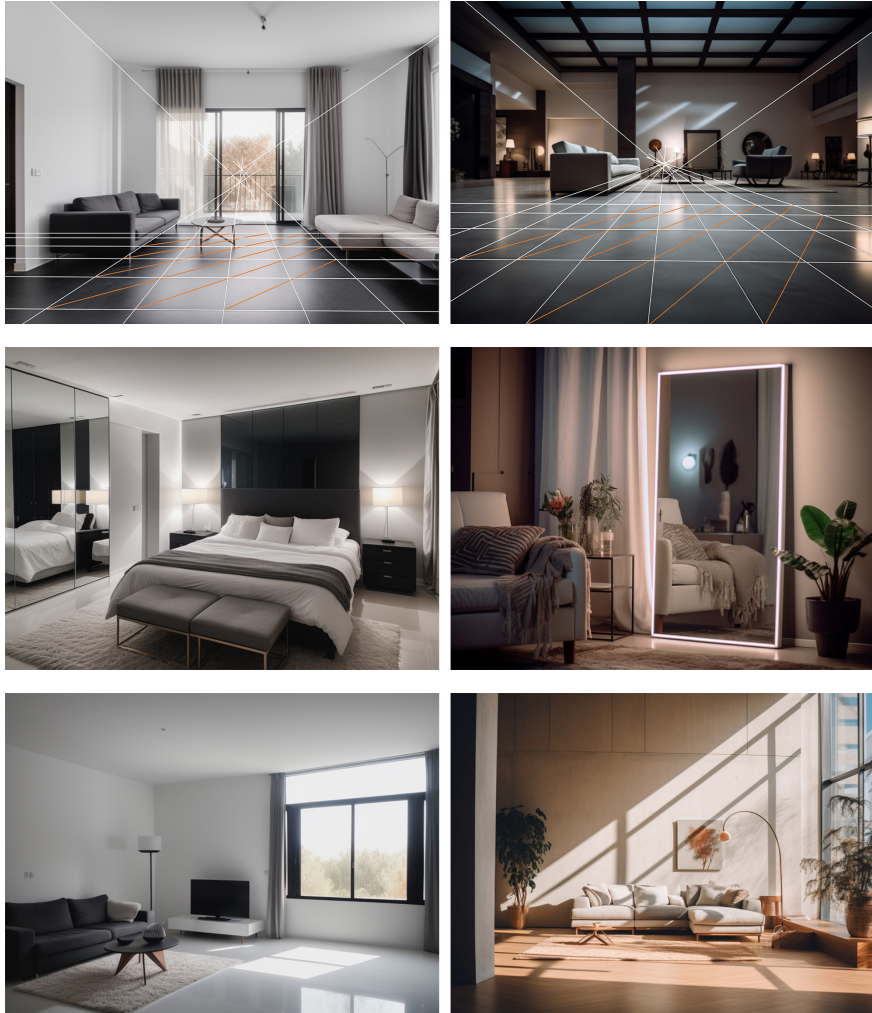


Fig. 10. Analisi prospettica di due delle precedenti immagini generate con StableDiffusion a sinistra e Midjourney a destra (elaborazione degli autori).

Fig. 11. Immagini generate con StableDiffusion a sinistra e Midjourney a destra attraverso il prompt: home interior view, bedroom, modern architecture, minimalist style furniture, mirror with reflections (elaborazione degli autori).

Fig. 12. Immagini generate con StableDiffusion a sinistra e Midjourney a destra attraverso il prompt: home interior view, modern architecture, minimalist style furniture, dramatic lighting and shadows (elaborazione degli autori).

L'analisi di elementi riflessi negli specchi ha presentato l'incoerenza di alcune soluzioni, sottolineando l'inconsapevolezza spaziale dell'AI. In particolare, il riflesso dello specchio manca di alcuni elementi presenti nella scena (come la coperta stesa sul letto in fig. 11 a sinistra) oppure mostra tali elementi in posizione incoerente (sempre la coperta in fig. 11 a destra, che nella scena è appoggiata al bracciolo del divano mentre nel riflesso è stesa dall'altro lato).

Lo studio delle ombre è un'ulteriore conferma del fatto che la costruzione delle immagini generate dall'AI non abbia consapevolezza dello spazio tridimensionale che rappresenta. Molto spesso i raggi di luce che penetrano dalle finestre producono delle ombre che non sono coerenti con gli infissi (fig. 12).

Conclusioni

L'analisi del funzionamento delle AI *text-to-image*, insieme alla ricognizione degli studi sul tema e alle sperimentazioni condotte, permette di tracciare delle prime riflessioni sul loro possibile ruolo per la rappresentazione del progetto di architettura.

Le sperimentazioni mettono in risalto sia potenzialità che punti di debolezza. Tra le prime possono certamente essere annoverate la grande velocità di generazione che permette, in poche decine di secondi, di disporre di immagini dall'elevata qualità visuale, la flessibilità delle tecniche grafiche (ri)prodotte e la coerenza con i *prompt* proposti. Queste potenzialità rendono possibili dei rapidi salti avanti-indietro nel processo progettuale, dalla fase preliminare ideativa a quella di visualizzazione avanzata dell'idea. Tra le seconde, oltre ai già analizzati limiti in termini di correttezza della rappresentazione, bisogna annoverare le criticità segnalate da più parti riguardo alla legittimità in termini di diritti d'autore dei metodi adottati per creare i database per l'allenamento [5] e la presenza di potenziali *bias* culturali indotti nelle AI. A tal proposito, è sufficiente notare come le piazze proposte in fig. 8 al centro rispecchino chiaramente modelli nord americani, lontani dalla concezione europea di spazio pubblico. Inoltre, si evidenzia come questo tipo di AI non siano

attualmente idonee per contribuire alla realizzazione di elaborati tecnici di progetto.

Esiste una sostanziale differenza concettuale tra l'intelligenza artificiale e quella umana, che consegue dal processo di allenamento e generazione delle AI. Queste ultime sono intelligenze di tipo interpolativo, cioè efficientissime nell'interpolare valori esistenti nell'ambito del database di allenamento e generare un valore non presente ma che non gli sarà mai del tutto estraneo. Non sono cioè in grado di estrapolare nuovi valori, non solo non presenti ma del tutto alieni al database. Questa forma di intelligenza è, invece, tipicamente umana [Del Campo 2022a]. La potenzialità della co-creazione uomo-AI in fase ideativa sembra quindi risiedere proprio nella collaborazione tra due tipi diversi di intelligenze, in cui quella interpolativa, avviata e guidata dagli input umani, propone immagini «familiari ma strane» [Del Campo 2022b, p. 28] dalle quali l'intelligenza umana può cogliere suggestioni da sviluppare in idee innovative. Questa ipotesi rinnova la problematica riguardante l'autorialità che già la progettazione/rappresentazione algoritmica ha aperto, portando all'idea che un'autorialità condivisa da più agenti (umani o artificiali) sia connaturata al progredire della rivoluzione digitale in architettura. In questo contesto, l'autorialità umana va comunque intesa come "primaria" poiché ha il ruolo di creare le regole generali, gli "oggetti" deleuziani, dalle quali le autorialità artificiali 'secondarie' deriveranno le singole forme, gli oggetti [Carpo 2011, p. 40, 123-128]. La transizione dalla figura di architetto come progettista di singole forme a quella di progettista di regole generali è già in corso e ha portato all'ampliamento dello spettro dei linguaggi adottati. A partire dalla seconda svolta digitale, gli architetti hanno imparato a comporre *script* e algoritmi, affiancandoli alla rappresentazione grafica, e ora, con l'avvento delle AI basate su *prompt*, sono chiamati a integrare il linguaggio naturale in forma scritta tra i loro metodi di progettazione.

Quest'ultima sfida posta dalla rivoluzione digitale deve far riflettere sui linguaggi di rappresentazione in senso più ampio e sul loro insegnamento, uno dei possibili campi di ricerca interdisciplinare del presente e del prossimo futuro del disegno d'architettura.

Crediti

Gli autori hanno condiviso tutte le fasi della ricerca in modo equo. Ai fini della stesura dell'articolo: M.F.M. ha scritto *Introduzione, Avanti e indietro*

nello spazio latente, Potenzialità e Conclusioni; S.M. ha scritto *Passato e presente, Le principali piattaforme, Studi relativi e Limiti*.

Note

[1] Dato relativo a luglio 2023, fonte <<https://discord.com/servers>> (consultato il 24 luglio 2023).

[2] <<https://rom1504.github.io/clip-retrieval/?back=https%3A%2F%2Fknn.laion.ai&index=laion5B-H-14&useMclip=false>> (consultato il 24 luglio 2023).

[3] La sperimentazione è stata condotta su due piattaforme di AI: *Midjour*

ney e *StableDiffusion* associato a un ulteriore modello di addestramento chiamato *RealisticVision* v. 4.

[4] Nel *prompt* è stata anche inserita una specifica relativa alla presenza di un pavimento composto da piastrelle quadrate in modo da permettere la successiva analisi prospettica.

[5] <<https://www.egaire.eu/>> (consultato il 24 luglio 2023).

Autori

Matteo Flavio Mancini, Dipartimento di Architettura, Università degli Studi Roma Tre, matteoflavio.mancini@uniroma3.it
Sofia Menconero, Dipartimento di Storia, Disegno e Restauro dell'Architettura, Sapienza Università di Roma, sofia.menconero@uniroma1.it

Riferimenti bibliografici

Carpó, M. (2011). *The alphabet and the algorithm*. Cambridge - London: The MIT Press.

Carpó, M. (ed.). (2013). *The digital turn in architecture 1992-2012*. Chichester: John Wiley & Sons.

Carpó, M. (2017). *The second digital turn: design beyond intelligence*. Cambridge - London: The MIT Press.

Colton, S. et al. (2021). Generative Search Engines: Initial Experiments. In A. Gómez de Silva Garza et al. (a cura di). *Proceedings of the 12th International Conference on Computational Creativity*, Mexico City, 14-18 settembre 2021, pp. 237-246. Mexico City: ACC.

Crowson, K. et al. (2022). *VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance*. In *arXiv*. <<https://arxiv.org/abs/2204.08583>> (consultato il 18 luglio 2023).

Del Campo, M. (2022a). When Robots Dreams. In *Conversation with Alexandra Carlson*. In *Architectural Design*, n. 03, v. 92, pp. 47-53.

Del Campo, M. (2022b). *Neural Architecture. Design and Artificial Intelligence*. Novato: Oro Editions.

Dhariwal, P., Nichol, A. (2021). Diffusion Models beat GANs on Image Synthesis. In M. Ranzato et al. (eds). *Advances in Neural Information Processing Systems*, v. 34, pp. 1-15. Cambridge: MIT Press.

Goodfellow, I. et al. (2014). Generative Adversarial Nets. In Z. Ghahramani et al. (eds). *Advances in Neural Information Processing Systems*, v. 29, pp. 1-9. Cambridge: MIT Press.

Hegazy, M., Saleh, A.M. (2023). Evolution of AI role in architectural design: from parametric exploration and machine hallucination. In *MSA Engineering Journal*, v. 2, n. 2, pp. 262-288. <www.doi.org/10.21608/MSAENG.2023.291873> (consultato il 18 luglio 2023).

Hong, W. et al. (2022). CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers. In *arXiv*. <<https://arxiv.org/abs/2205.15868>> (consultato il 18 luglio 2023).

Jaruga-Rozdolska, A. (2022). Artificial intelligence as part of future practices in the architect's work: Midjourney generative tool as part of a process of creating an architectural form. In *Architectus*, v. 3, n. 71, pp. 95-104.

Jun, H., Nichol, A. (2023). *Shape-E: Generating Conditional 3D Implicit Functions*. In *arXiv*. <<https://arxiv.org/abs/2305.02463>> (consultato il 18 luglio 2023).

Nichol, A. et al. (2022). *Point-E: A System for Generating 3D Point Clouds from Complex Prompts*. In *arXiv*. <<https://arxiv.org/abs/2212.08751>> (consultato il 18 luglio 2023).

Paananen, V. et al. (2023). Using Text-to-Image Generation for Architectural Design Ideation. In *arXiv*. <<https://arxiv.org/abs/2304.10182>> (consultato il 18 luglio 2023).

Ploennings, J., Berger, M. (2022). AI Art in Architecture. In *arXiv*. <<https://arxiv.org/abs/2212.09399>> (consultato il 18 luglio 2023).

Prix, W. et al. (2022). The Legacy Sketch Machine. From Artificial to Architectural Intelligence. In *AD, Machine Hallucinations: Architecture and Artificial Intelligence*, v. 92, n. 3, pp. 14-21.

Ramesh, A. et al. (5 gennaio 2021). DALL-E: Creating images from text. <<https://openai.com/research/dall-e>> (consultato il 18 luglio 2023).

Radford, A. et al. (2021). Learning Transferable Visual Models from Natural Language Supervision. In M. Meila, T. Zhang (a cura di). *Proceedings of the 38th International Conference on Machine Learning*. Virtuale, 18-24 luglio, v. 139, pp. 8748-8763. Maastricht: ML Research Press.

Reed, S. et al. (2016). Generative Adversarial Text to Image Synthesis. In M. F. Balcan, K. O. Weinberger (a cura di). *Proceedings of the 33rd International Conference on Machine Learning*, v. 48, pp. 1060-1069. Maastricht: ML Research Press.

Rombach, R. et al. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, New Orleans, 18-24 giugno, pp. 10674-10685. New York: IEEE.



Singer, U. et al. (2022). *Make-A-Video: Text-to-Video Generation without Text-Video data*. In *arXiv*. <<https://arxiv.org/abs/2209.14792>> (consultato il 18 luglio 2023).

Tong, H. et al. (2023). An attempt to integrate AI-based techniques into first year design representation course. In K.Vaes, J.Verlinden (a cura di). *Connectivity and Creativity in times of Conflicts. Cumulus Conference Proceedings*. Anversa, 12-15 aprile, pp. 1-5. Anversa: University of Antwerp.

Tsigkari, M. et al. (29 marzo 2021). Towards Artificial Intelligence in Architecture: How machine learning can change the way we approach design. In *Plus Journal*, <<https://www.fosterandpartners.com/insights/plus-journal/towards-artificial-intelligence-in-architecture-how-machi->

[ne-learning-can-change-the-way-we-approach-design](https://www.fosterandpartners.com/insights/plus-journal/towards-artificial-intelligence-in-architecture-how-machine-learning-can-change-the-way-we-approach-design)> (consultato il 18 luglio 2023).

Wallish, S. (2022). GAN Hadid. In S. Carta (a cura di). *Machine Learning and the City: Applications in Architecture and Urban Design*, pp. 477-481. Hoboken-Chichester: John Wiley & Sons.

Yildirim, E. (2022). Text-to-image generation A.I. in architecture. In H. Hale Kozlu (a cura di). *Art and Architecture: Theory, Practice and Experience*, pp. 97-119. Lyon: Livre de Lyon.

Zhang, L., Agrawala, M. (2023). Adding Conditional Control to Text-to-Image Diffusion Models. <<https://arxiv.org/abs/2302.05543>> (consultato il 18 luglio 2023).

AI-aided Design? Text-to-image Processes for Architectural Design

Matteo Flavio Mancini, Sofia Menconero

Abstract

Artificial Intelligence (AI) is marking a turning point in many aspects of human life, and it is appropriate to question its potential use in the architectural representation processes. This contribution provides a brief overview of the recent past of AI technologies to explain how they work, a snapshot of the current state of the art from text-to-image processes to image-to-3D processes, mainly focusing on the StableDiffusion platform. It also offers an overview of the latest studies in the field of architectural design. The subsequent experimentation becomes an opportunity to showcase the potential of AI in the co-creation process and the ability to simulate various graphic techniques, up to photorealistic visualization. On the other hand, it presents the limitations that, at the current stage of development, sometimes invalidate the results of text-to-image processes concerning the scientific aspects of representation. The conclusions reflect on the differences between human and artificial intelligence, the theme of shared authorship between humans and machines, and their consequences for architectural design.

Keywords: artificial intelligence, text-to-image, design drawing, authorship, stablediffusion.

Introduction

Architecture and architectural drawing have undergone significant developments over the past thirty years. The *first digital turn* [Carpo 2013] introduced digital representation in the 1990s, while the *second digital turn* [Carpo 2017] began with the spread of algorithms and big data in the 2010s. Ten years later, we are witnessing another potential turning point due to the sudden development and diffusion of Artificial Intelligence (AI) tools, which are already in use in major architectural firms such as Coop Himmelb(l)au [Prix et al. 2022], Zaha Hadid Architects [Wallish 2022], and Foster + Partners [Tsigkari et al. 2021]. One branch of AI, based on text-to-image processes, offers easily accessible solutions dedicated to image

creation. This machine-learning model uses descriptive natural language as input and produces an image based on the provided description. The results obtained from these platforms are remarkable in matching the entered textual prompts and the flexibility of graphic techniques they can (re)produce. Starting from the assumption that, at the current state, these AIs have no creative consciousness or actual ability to understand compositional and projective rules or the spatiality represented in the images, it is still worth exploring their potential use in architectural representation processes. With this goal in mind and considering the intrinsic characteristics of this technology, which will be explained in the follow-

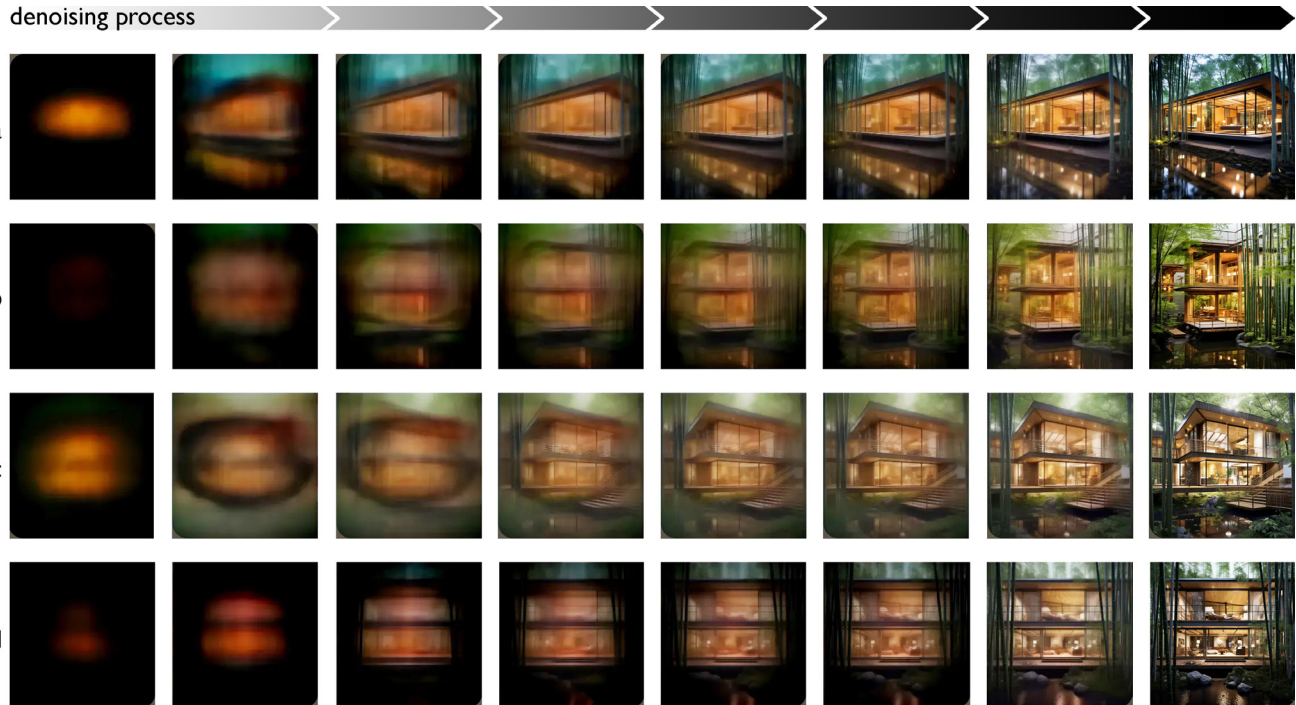
ing paragraphs, the experimentation of text-to-image AI through the open-source platform *StableDiffusion* is proposed to create perspective images capable of contributing to the preliminary stages of project conception.

(Recent) past and present of AI-based image generation

Generative AI systems in the field of architecture and design are rapidly advancing. A pivotal moment in image generation research was marked by the invention of Generative Adversarial Networks (GANs) in 2014 [Goodfellow et al. 2014]. This represents a deep learning architecture where two opposing neural networks, a generator and a discriminator, inter-

act iteratively during training to reach a point where the discriminator can no longer differentiate between synthetic images produced by the generator and real images used as training data. In 2016, a GAN architecture capable of generating plausible images from detailed textual descriptions was developed [Reed et al. 2016], effectively initiating the AI text-to-image system. Another significant advancement came with a more efficient natural language processing-based learning method known as CLIP (Contrastive Language-Image Pre-training) [Radford et al. 2021]. This image classification model identifies objects by learning from text associated with an image, rather than relying on manually assigned labels, and it was trained on 400 million image-text pairs extracted from the web. CLIP models can estimate the conformity of a generated image to a textual prompt [Colton et al. 2021].

Fig. 1. Denoising process during the generation of four variants (a, b, c, d) in Midjourney through the prompt: a modern Japanese house in a bamboo forest in spring (authors' processing).



An example of this pairing is the image generation system called *VQGAN-CLIP*, which utilizes an even more powerful GAN neural network. The significant contributions of the *VQGAN-CLIP* architecture include visual quality in both image generation and manipulation, semantic fidelity between text and the generated image, efficiency due to not requiring additional training beyond pre-trained models, and the value of open development and scientific progress [Crowson et al. 2022, p. 2]. Subsequently, GAN-based systems were replaced by diffusion models, which are probabilistic machine learning models trained to remove noise previously introduced from images by learning to reverse the diffusion process [Dhariwal et al. 2021]. The training of these models enables them to utilize denoising methods to synthesize new noise-free images from random inputs (fig. 1).

Some applications of AI in the field of architecture and design include:

- Text-to-image: the most common operation involves generating images based on textual descriptions, often combined with other functionalities.
- Image-to-image: transforming an input image to match the characteristics of a target image. This can be used for style transfer, object manipulation (inpainting), converting black and white images to color, or increasing image resolution (upscaling).
- Text or image-to-video: creating videos from textual prompts (e.g., *Make-a-Video* [Singer et al. 2022] or *CogVideo* [Hong et al. 2022]) or generating animations by editing images generated through image-to-image (e.g., *Deforum*), producing effects similar to stop-motion videos.
- Text or image-to-3D: generating 3D models from textual prompts (e.g., *Point-E* for point clouds [Nichol et al. 2022] or *Shape-E* for textured meshes [Jun et al. 2023]), or generating 3D models from images (e.g., *Kaedim*).

The recent incredible proliferation of AI text-to-image capabilities can be attributed to the activation of user-friendly platforms, even for non-expert users, such as *DALL-E 2*, *Midjourney*, and *StableDiffusion*.

The main platforms for text-to-image AI

DALL-E is the first of the three platforms introduced

in January 2021 (the current version 2 was released in April 2022) by *OpenAI* [Ramesh et al. 2021], the same developers of *ChatGPT*. The platform, available for online subscription, offers four functions: the generation of realistic and artistic images from a textual description that can combine concepts, attributes, and styles (fig. 2); outpainting, which involves extending the image beyond its original boundaries by creating a new composition; inpainting, which allows the modification of image portions by adding or removing objects through a textual description while maintaining consistency with the rest of the scene; and the generation of variations inspired by an input image.

Midjourney, released on July 12, 2022, has now reached version 5.2 with significant improvements in terms of adherence to prompts and photorealism compared to its initial release (fig. 3). After a year, it boasts over 15 million users [1]. Like *DALL-E*, it is available for online subscription. The three main generative activities on *Midjourney* are: generating an image from a textual prompt, generating a description from an image, and

Fig. 2. Image generated with *DALL-E 2* through the prompt: a modern building on a crowded street at sunset (authors' processing).



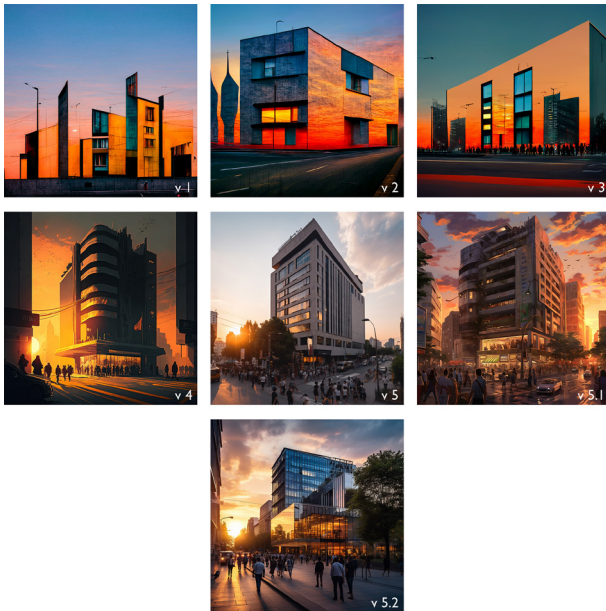


Fig. 3. Comparison between different images generated from the same text (prompt: a modern building on a crowded street at sunset) for different versions of Midjourney (authors' processing).

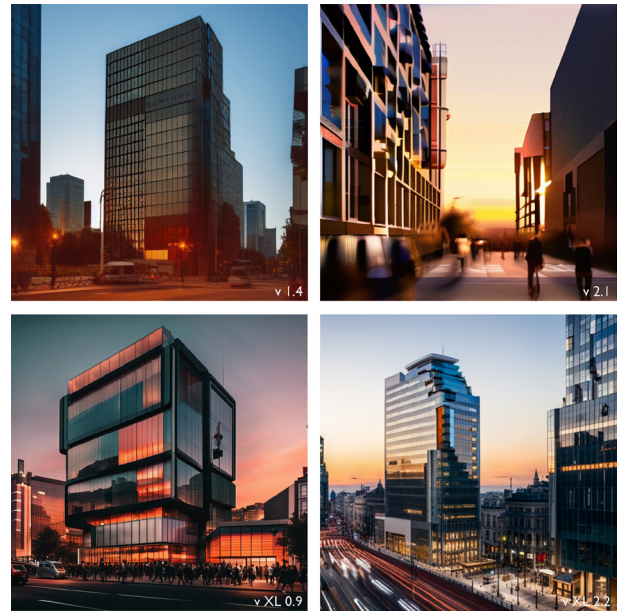


Fig. 4. Comparison between different images generated from the same text (prompt: a modern building on a crowded street at sunset) for different versions of StableDiffusion (authors' processing).

synthesizing an image from two to five input images. *Midjourney* (like *StableDiffusion*) also allows the use of a negative prompt in case specific elements are not desired in the generated image.

StableDiffusion, released in August 2022, is the only one among the three platforms to be open-source and is based on a diffusion model called *latent diffusion model* [Rombach et al. 2022]. The current beta XL version is only available online through a subscription, while previous versions can be installed locally for free. *StableDiffusion* supports image generation using a text prompt describing the elements to include or exclude from the output (fig. 4), inpainting and outpainting, image-to-image generation, and upscaling. It is also possible to add extensions to *StableDiffusion*, such as *ControlNet*, which generates variations of an input image through textual descriptions, and *Deforum*, which, through the image-to-image function, generates a series of images with minor transformations and stitches

them together to create a video.

The biggest advantage of *StableDiffusion* over the other platforms is that end-users can implement additional training (fine-tuning) to optimize generation outputs to specific use cases. For example, in architectural studies where AI has become part of the creative process, the neural network is trained with targeted images from the studio's design repertoire to achieve results more in line with architectural and graphic language.

Therefore, unlike the previous platforms, *StableDiffusion* allows greater freedom in terms of customization of the generative process, which is why it was chosen for further experimentation, coupled with the *ControlNet* extension [Zhang et al. 2023], which improves output control. The latter is a neural network structure designed to handle diffusion models by incorporating additional conditions: by manipulating the input conditions of the blocks, it further controls the overall

behavior of an entire neural network. *ControlNet* operates based on an input image and a textual description, allowing the generation of images that conform compositionally to the input but also follow the specified description. The data processing involves, first and foremost, the generation of a map based on the input image (called the annotation or preprocessing phase), which is used by the network to generate variants with the described textual characteristics (fig. 5).

Back and forward in latent space: between training and generation in the *StableDiffusion* model

To properly approach the experiments that will follow in the next paragraphs, it is essential to try to understand not so much the strictly technical-computer aspects but the processes carried out by this type of AI, specifically *StableDiffusion*, in the two distinct moments of training and generation, as both appropriate use and critical interpretation of this technology depend on them.

Diffusion models borrow from thermodynamics the concept of diffusion, which is the phenomenon whereby particles of a fluid move randomly within another fluid with a different concentration until they reach a new equilibrium condition. Similarly, AI-generated images during generation progressively emerge from the chaos of digital noise. The principle of diffusion is used both in the training phase (forward diffusion) and in the generation phase (reverse diffusion). In *StableDiffusion*, both processes occur in the latent space, a numerical/informational space in which images are translated into tensors (multi-dimensional matrices) to work on a compressed version of them, lighter than the initial pixel space of the images. Texts describing the images also undergo a similar translation and compression. The analogy between the latent representations of texts and images is important because it helps understand that AIs do not store and process collections of syllables, words, or portions of images but operate on abstract numerical representations of image features, represented objects, possible environmental situations, and various techniques and styles. The

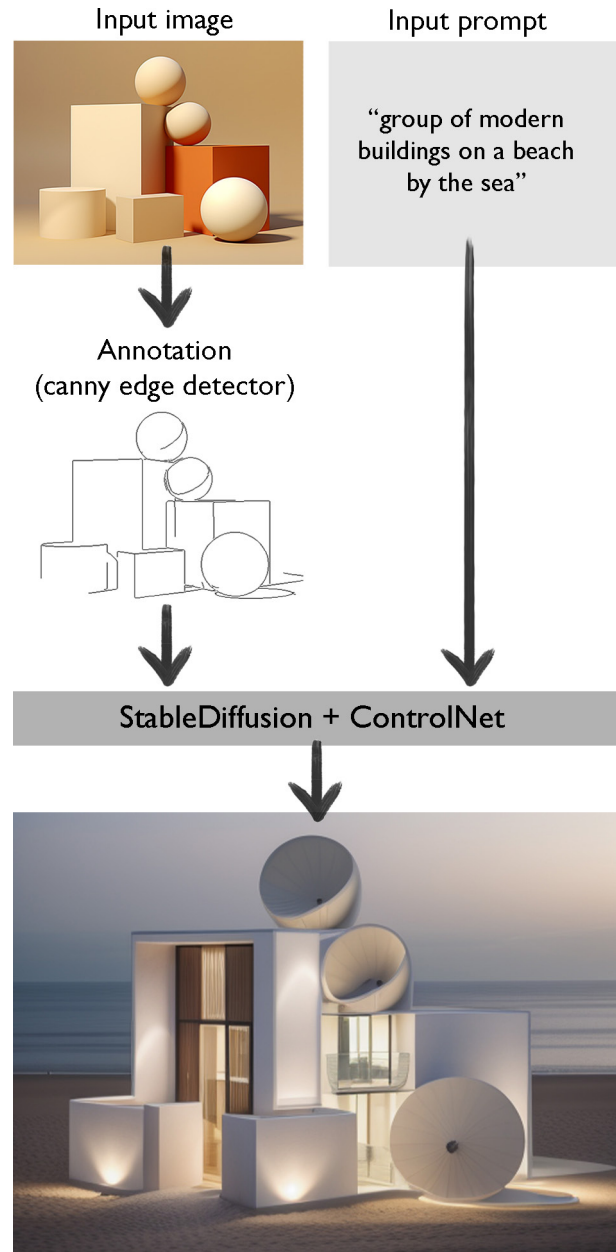


Fig. 5. Diagram of image generation through *StableDiffusion* with *ControlNet* conditioning (canny edge) (authors' processing).

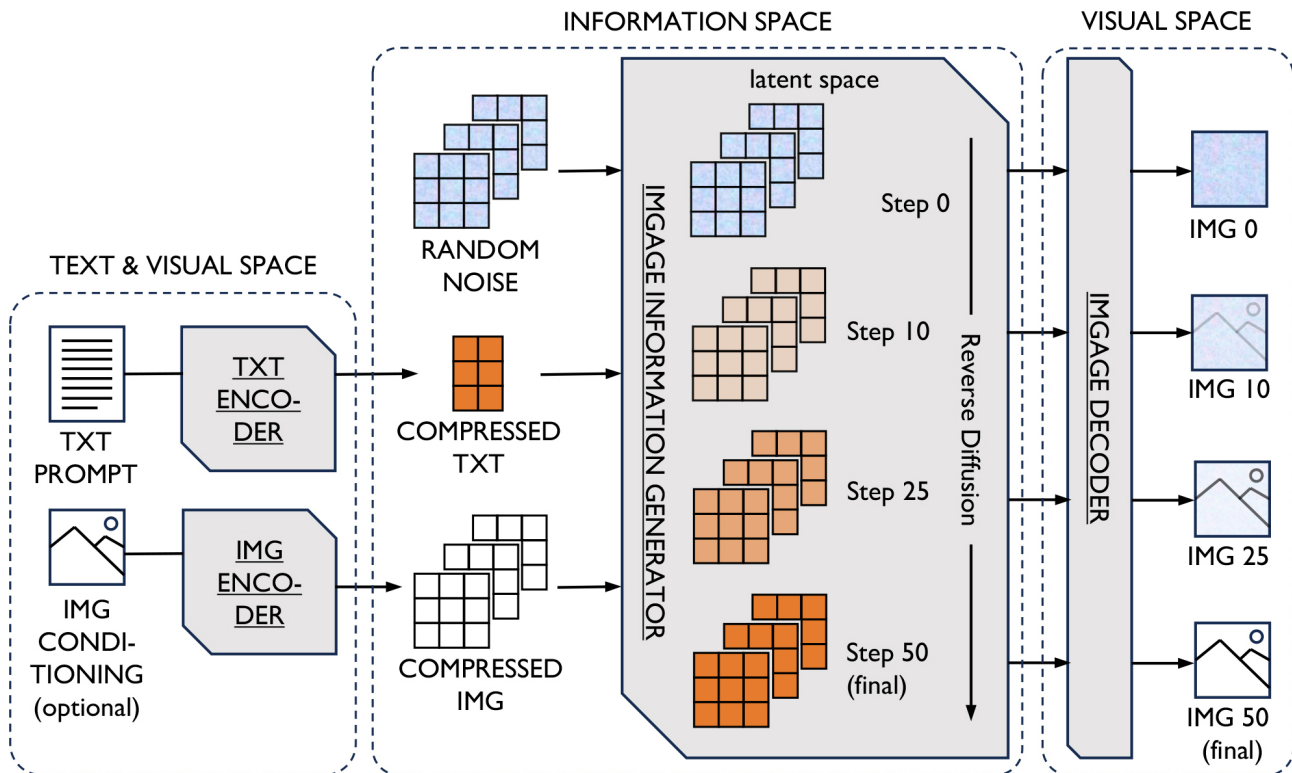
latent space can be imagined as the place where the AI stores, during training, and retrieves, during generation, its 'knowledge.'

The training of these AIs is not progressive over time but occurs before publication, so their 'knowledge' is static and periodically updated with the progression of versions. For example, *StableDiffusion* is trained on the open-access dataset LAION 5B, composed of 5 billion image-caption pairs whose content is explorable, starting from the input of a textual prompt, through a dedicated portal [2]. Consulting the training dataset allows one to get an idea of the correspondences between terms and images and, therefore, what can be expected from the results of the processing: a search that does not return consistent results indicates that

images corresponding to the expectations for that textual input cannot be generated. During training, images from the database are processed by introducing random noise patterns of different intensities. The images thus processed, together with the corresponding captions, are subjected to the AI to train it to identify the type of pattern adopted, the amount of noise introduced and to remove both to improve the quality of the images. In this way, through forward diffusion, the AI simultaneously learns how to obtain noise-free images and the correspondences between images and texts.

What is learned is used by *StableDiffusion* to develop a generative process that starts and ends in a space where data (texts and images) are suitable for human

Fig. 6. Diagram of the image generation process through the diffusion model adopted by *StableDiffusion* (authors' processing).



perception, passes through a purely informational space (latent space) where data are represented by tokens (texts) and tensors (images).

The generative process can be divided into three fundamental blocks (fig. 6). The first involves the compression and transformation into numerical information through encoders (specifically trained neural networks) of the input data entered to condition image generation. In *StableDiffusion*, thanks to the *ControlNet* extension, textual inputs can be supplemented with optional graphical conditions. In the second block, through the reverse diffusion process, the processing of the inputs in relation to known knowledge takes place. This iterative process goes through several denoising steps to refine the correspondence between the input entered and the generated image. At this stage, processing occurs at the level of numerical information, and there is no graphic image processing. The latter occurs in the third block, where numerical representations are translated from a neural network with a decoding function into visually perceptible images [Rombach et al. 2022].

Related studies on AI applied to architectural design

Some studies related to AI in the field of architectural design are focused on highlighting the potential and limitations of the technology. In most cases, the potential is recognized as a support in the creative process [Jaruga-Rozdolska 2022; Paananen et al. 2023]. Among other potential uses, the ability to imagine abstract forms, reimagine biomimetic architecture, revisit traditional architecture, and visualize photo-realistic advancements starting from architectural sketches are mentioned. The identified limitations are related to the possibility of control and customization of processes, insufficient consideration of structural feasibility aspects, and potential stylistic-architectural inconsistency in the generated results [Hegazy et al. 2023]. Case studies related to architectural projects where AI has been applied include the ideation phase, the generation of sketches with specific graphic styles, the addition of people and objects to existing images, the combination of various parts of images into a coherent composition, the variation of an initial image, the change of the graphic style of an existing

image, floor plan design, exterior and interior design, texture creation, and urban planning [Ploenning et al. 2022; Yildirim 2022].

One educational study involves integrating AI techniques with traditional techniques in a first-year university design representation course, where the authors observed an improvement in students' interpretative and compositional abilities [Tong et al. 2023]. Students were asked to create a composition of solids and draw hand-drawn orthogonal projections and isometric axonometry. Then, they were required to generate a series of images using *Midjourney* with specific keywords. Finally, they were instructed to combine the two previous productions using various techniques.

Potential of text-to-image AI for architectural drawing

To experiment with the potential contribution of AI in the preliminary phase of the project, the moment when representation contributes to ideation and prefiguration, it was decided to work on idea definition and visualization.

Three different graphic inputs were hypothesized: two external perspective views of a three-dimensional volumetric model and a sketch of an interior, intentionally lacking characterizations except for the minimum necessary for spatial definition and framing. These graphic inputs, thanks to the *ControlNet* extension, are entrusted with incorporating the general morphological setting of the project into the generative process, while textual inputs are used to describe the desired graphic techniques and any architectural features related to materials, context, and additional stylistic characteristics that one wishes to include. The results of these initial experiments demonstrate the remarkable flexibility of AI in (re)creating different graphic techniques, ranging from pencil drawings to colored pencils to watercolors, with a significant ability to integrate both natural and artificial context elements. Simultaneously, the AI's addition of detailed elements such as textures, perforations, and materials contributes to the advancement of ideation in a process where it can be hypothesized that some of these elements may actually be inserted into the

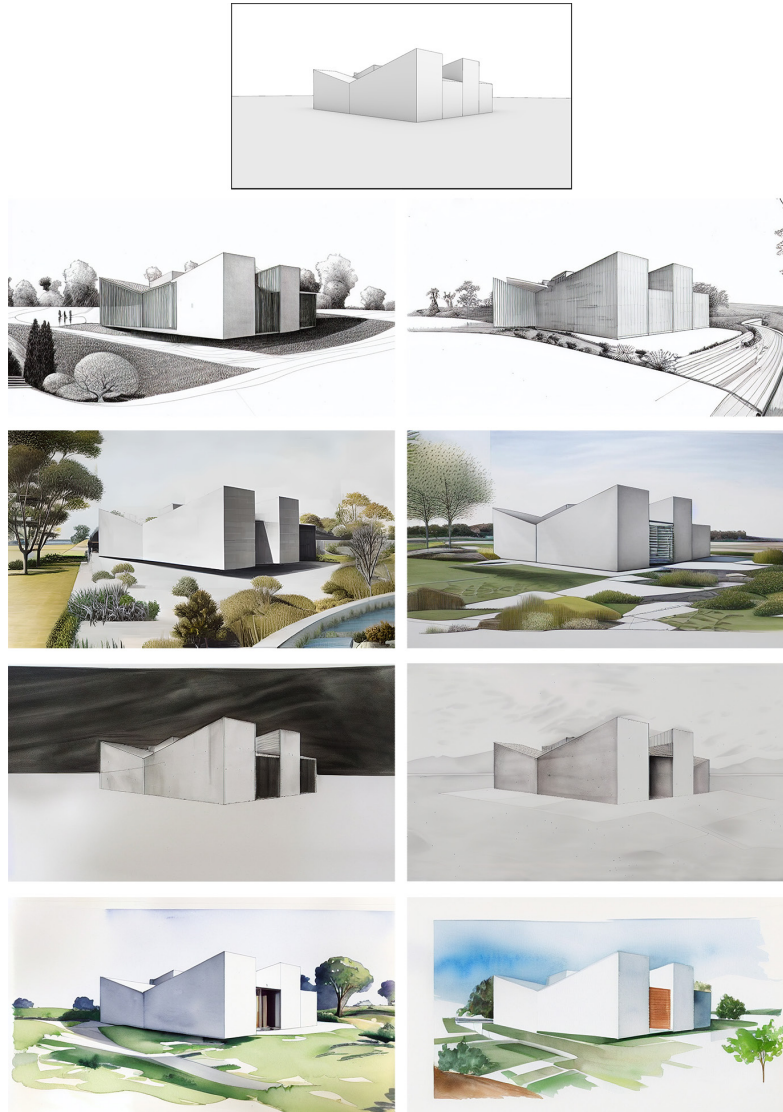


Fig. 7. Images generated with StableDiffusion to simulate various graphic techniques. From top to bottom: pencil, colored pencils, monochromatic watercolor, and colored watercolor. Prompt: linear, exterior view, contemporary architecture, highly detailed architecture, large windows, concrete, architectural drawings, technical drawings, [desired graphic technique], line drawings, working drawings, architectural sketches, conceptual style, abstract (authors' processing).

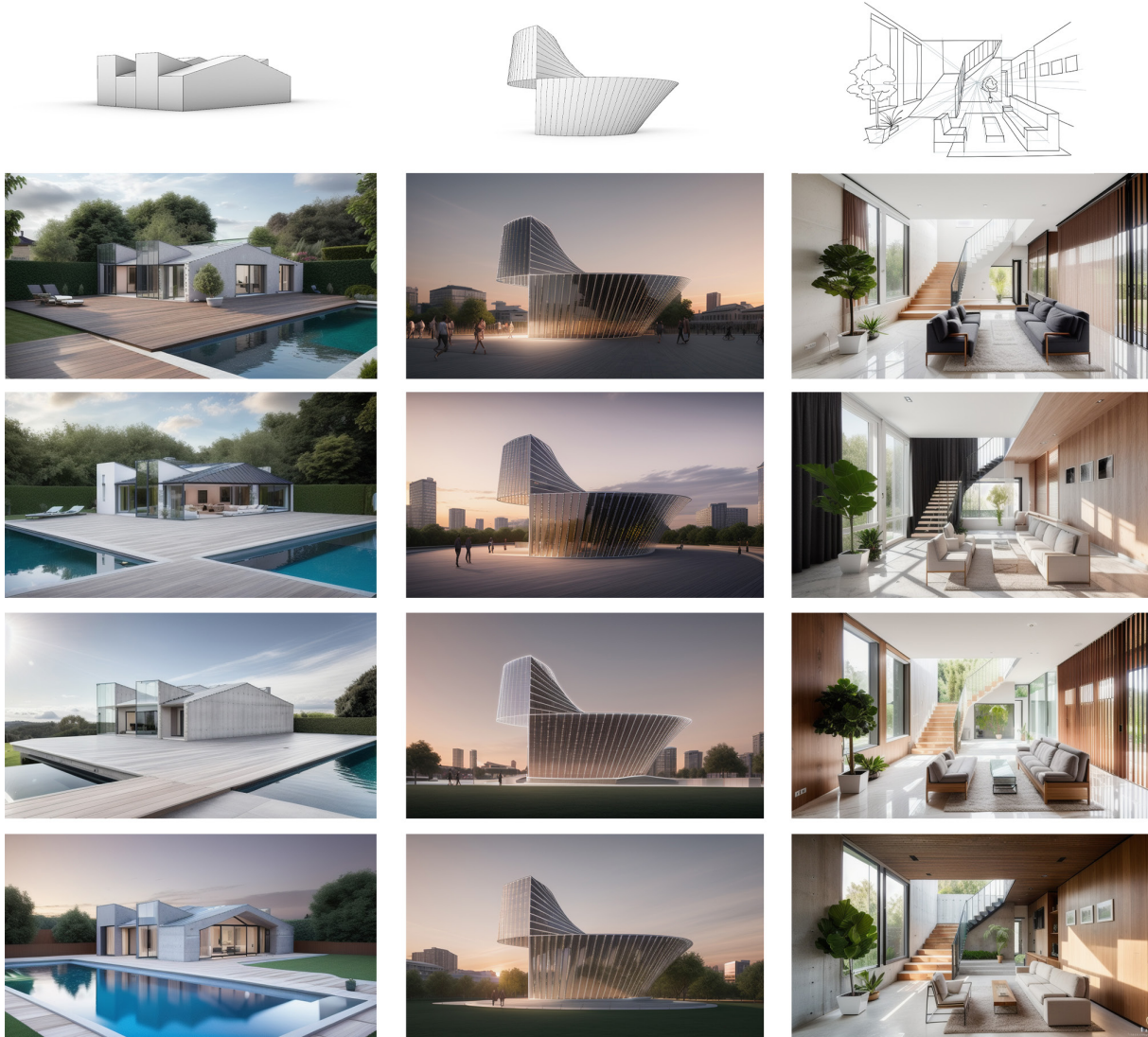
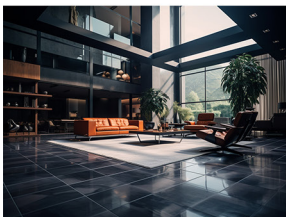
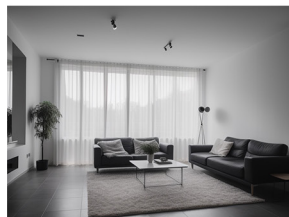


Fig. 8. Images generated with StableDiffusion for photorealistic visualizations. On the left, exterior views, prompt: exterior home view, concrete walls and roof, large glass windows, small rectangular swimming pool, garden, garden furniture, clouds. In the center, exterior views with curved surfaces, prompt: a pavilion in a contemporary architecture style, covered with reflective panels, in a square with people and trees. On the right, interior views, prompt: home interior view, modern architecture, large glass windows with curtains, timber framing, wood flooring, concrete ceiling, steel staircase, large sofa with pillows, armchair, coffee table with flowerpot, carpet, plants, lamp, minimalist style furniture, sunlight from windows, daylight (authors' processing).

Midjourney



StableDiffusion + RealisticVision v4



ongoing project in a repeated human-machine exchange (fig. 7).

The potential contribution in terms of idea definition through the fast generation of variations is more evident when requesting the AI to produce photorealistic images. In this case, the ability to propose variations based on the textual prompt is more recognizable. Experiments conducted on external views demonstrate the variety of materials and interpretations of the simple volumetric schemes provided as input, as well as the skill in creating contextual settings (fig. 8). Similarly, experiments based on a digital sketch of an interior environment highlight the ability to combine colors and materials but also the inclination to add elements such as curtains and floor elevations. Minor elements such as light points and furnishings also appear. The distribution of these integrations generally aligns with the overall setting.

Limitations of text-to-image AI in representation

The limitations investigated in this paragraph [3] particularly concern the aspect of architectural representation (perspective, reflections, lighting/shadows). At present, AI has no awareness of the projective rules underlying correct perspective construction. While this assertion could be deduced based on the theoretical principles behind the technology, it is also experimentally confirmed. By adding a descriptive part about the representation method (central perspective) [4] in the textual prompt, we arrive at results where central perspective is present only in some of the generated images (fig. 9). When analyzing the perspective setting of two of the previous generated images, it is observed that the vanishing lines of the square floor tiles (horizontal lines perpendicular to the picture plane) do not converge to a unique point (fig. 10). Additionally, by tracing diagonals from the two visible ends of the square tiles in the images, it can be noted that the intermediate intersections are not perfectly aligned with the diagonals. Therefore, the perspectives are perceptually effective but not correct as projections. The results of the perspective experimentation suggest that the AI has not been trained to correctly recognize different representation methods.

Fig. 9. Images generated through the prompt: central perspective, home interior view, floor with regular dark square tiles, modern architecture, minimalist style furniture, daylight (authors' processing).

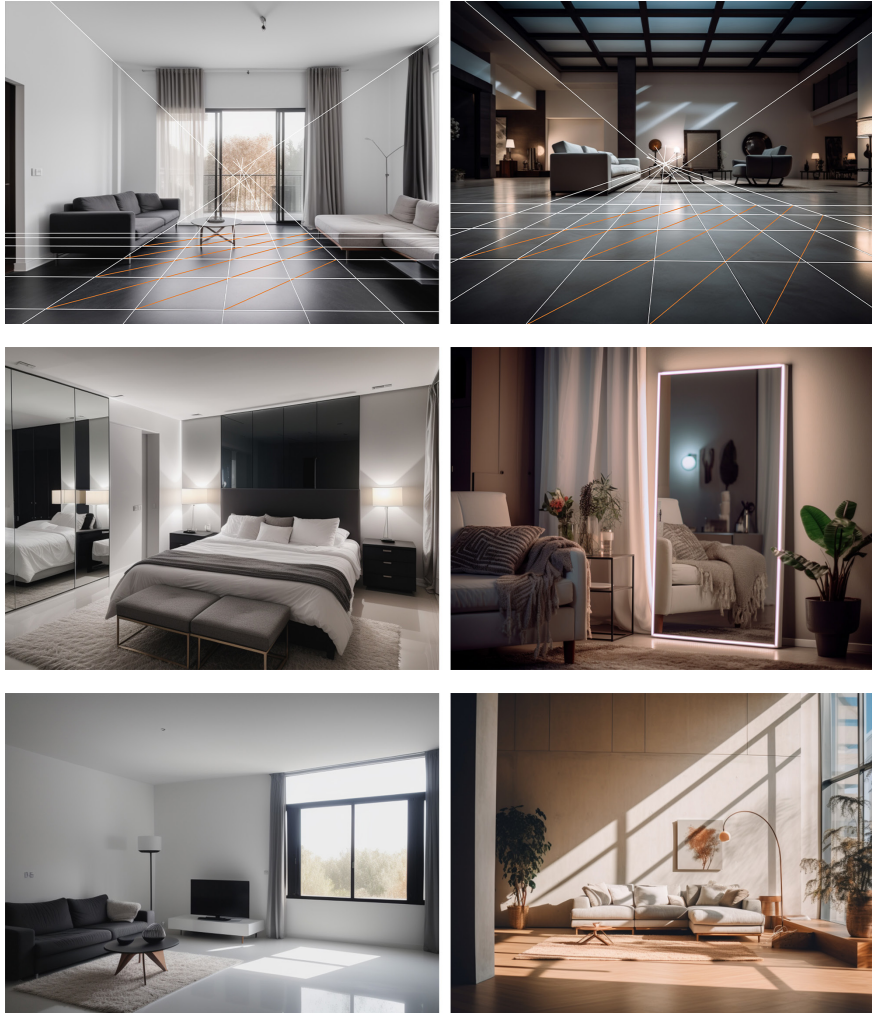


Fig. 10. Perspective analysis of two of the previous images generated with StableDiffusion on the left and Midjourney on the right (authors' processing).

Fig. 11. Images generated with StableDiffusion on the left and Midjourney on the right through the prompt: home interior view, bedroom, modern architecture, minimalist style furniture, mirror with reflections (authors' processing).

Fig. 12. Images generated with StableDiffusion on the left and Midjourney on the right through the prompt: home interior view, modern architecture, minimalist style furniture, dramatic lighting and shadows (authors' processing).

The analysis of elements reflected in mirrors revealed inconsistencies in some solutions, highlighting the spatial unawareness of AI. In particular, the mirror reflection lacks some elements present in the scene (such as the blanket on the bed in figure 11 on the left) or shows these elements in an inconsistent position (the same blanket in figure 11 on the right, which in the scene is draped over the sofa armrest, while in the reflection, it is spread out on the other side).

The study of shadows further confirms that the construction of images generated by AI lacks awareness of the three-dimensional space it represents. Very often, the light rays entering from windows produce shadows that are not consistent with the window fixtures (fig. 12).

Conclusions

The analysis of how text-to-image AIs function, the review of studies on the subject, and the experiments conducted allow us to outline initial reflections on AIs' possible role in architectural project representation. The experiments highlight both potentials and weaknesses. Among the strengths, one can certainly include the high generation speed, which allows for high-quality visual images to be produced in just a few seconds, the flexibility of (re)producing graphic techniques, and the coherence with the provided prompts. These potentials make it possible for rapid shifts back and forth in the design process, from the initial ideation phase to the advanced visualization of the idea. Among the weaknesses, in addition to the previously discussed limitations in terms of representation accuracy, there are concerns raised by multiple parties regarding the legitimacy of the copyright of the methods used to create the training databases [5] and the presence of potential cultural biases induced in the AIs. For instance, it is worth noting how the squares proposed in figure 8 reflect North American models, diverging from the European conception of public space. Furthermore, this AI type is unsuitable

Credits

The authors have equally shared all phases of the research. For the purpose of drafting the article: M.F.M. wrote the *Introduction, Back and Fore-*

ward in *Latent Space, Potential, and Conclusions*; S.M. wrote *(Recent) Past and Present, The Main Platforms, Related Studies, and Limitations*.

for creating technical project documents. There is a substantial conceptual difference between artificial intelligence and human intelligence, stemming from the training and generation processes of AIs. The latter is interpolative intelligence, highly efficient at interpolating existing values within the training database and generating a value that is not present but is never entirely alien to it. They cannot extrapolate new values, not only those not present but entirely foreign to the database. This second form of intelligence is typically human [Del Campo 2022a]. The potential of human-AI co-creation in the ideation phase seems to lie precisely in the collaboration between two different types of intelligences, in which the interpolative one, initiated and guided by human inputs, proposes "Familiar but Strange" images [Del Campo 2022b, p. 28] from which human intelligence can draw inspiration for developing innovative ideas. This hypothesis renews the issue of authorship that algorithmic design/representation has already raised, leading to the idea that shared authorship among multiple agents (human or artificial) is inherent in the progress of the digital revolution in architecture. In this context, human authorship should still be understood as 'primary' because it plays the role of creating general rules, the Deleuzian 'objectiles', from which 'secondary' artificial authorships will derive individual forms, the 'objects' [Carpo 2011, p. 40, 123-128]. The transition from the architect as the designer of individual forms to that of a designer of general rules is already underway and has expanded the spectrum of adopted languages. Starting from the *second digital turn*, architects have learned to compose scripts and algorithms alongside graphic representation, and now, with the advent of AI based on prompts, they are called upon to integrate written natural language into their design methods.

This latest challenge posed by the digital revolution should stimulate reflection on representation languages in a broader sense and their teaching, which is one of the possible fields of interdisciplinary research in the present and near future of architectural drawing.

Notes

[1] Data as of July 2023, source <<https://discord.com/servers>> (accessed 24 July 2023).

[2] <<https://rom1504.github.io/clip-retrieval/?back=https%3A%2F%2Fknn.laion.ai&index=laion5B-H-14&useMclip=false>> (accessed 24 July 2023).

[3] The experimentation was conducted on two AI platforms: *Mi-*

djourney and *StableDiffusion*, associated with an additional training model called *RealisticVision v. 4*.

[4] The prompt also included a specification regarding the presence of a floor composed of square tiles to enable subsequent perspective analysis.

[5] <<https://www.egair.eu/>> (accessed 24 July 2023).

Authors

Matteo Flavio Mancini, Department of Architecture, Roma Tre University, matteoflavio.mancini@uniroma3.it

Sofia Menconero, Department of History, Representation and Restoration of Architecture, Sapienza University of Rome, sofia.menconero@uniroma1.it

Reference List

Carpó, M. (2011). *The alphabet and the algorithm*. Cambridge - London: The MIT Press.

Carpó, M. (ed.). (2013). *The digital turn in architecture 1992-2012*. Chichester: John Wiley & Sons.

Carpó, M. (2017). *The second digital turn: design beyond intelligence*. Cambridge - London: The MIT Press.

Colton, S. et al. (2021). Generative Search Engines: Initial Experiments. In A. Gómez de Silva Garza et al. (Eds.). *Proceedings of the 12th International Conference on Computational Creativity*, Mexico City, 14-18 September 2021, pp. 237-246. Mexico City: ACC.

Crowson, K. et al. (2022). *VQGAN-CLIP: Open Domain Image Generation and Editing with Natural Language Guidance*. In *arXiv*. <<https://arxiv.org/abs/2204.08583>> (accessed 18 July 2023).

Del Campo, M. (2022a). When Robots Dreams. In *Conversation with Alexandra Carlson*. In *Architectural Design*, No. 03, V. 92, pp. 47-53.

Del Campo, M. (2022b). *Neural Architecture. Design and Artificial Intelligence*. Novato: Oro Editions.

Dhariwal, P., Nichol, A. (2021). Diffusion Models beat GANs on Image Synthesis. In M. Ranzato et al. (eds.). *Advances in Neural Information Processing Systems*, V. 34, pp. 1-15. Cambridge: MIT Press.

Goodfellow, I. et al. (2014). Generative Adversarial Nets. In Z. Ghahramani et al. (eds.). *Advances in Neural Information Processing Systems*, v. 29, pp. 1-9. Cambridge: MIT Press.

Hegazy, M., Saleh, A.M. (2023). Evolution of AI role in architectural design: between parametric exploration and machine hallucination. In *MSA Engineering Journal*, V. 2, No. 2, pp. 262-288. <www.doi.org/10.21608/MSAENG.2023.291873> (accessed 18 July 2023).

Hong, W. et al. (2022). *CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers*. In *arXiv*. <<https://arxiv.org/abs/2205.15868>> (accessed 18 July 2023).

Jaruga-Rozdolska, A. (2022). Artificial intelligence as part of future practices in the architect's work: MidJourney generative tool as part of a process of creating an architectural form. In *Architectus*, V. 3, No. 71, pp. 95-104.

Jun, H., Nichol, A. (2023). *Shape-E: Generating Conditional 3D Implicit Functions*. In *arXiv*. <<https://arxiv.org/abs/2305.02463>> (accessed 18 July 2023).

Nichol, A. et al. (2022). *Point-E: A System for Generating 3D Point Clouds from Complex Prompts*. In *arXiv*. <<https://arxiv.org/abs/2212.08751>> (accessed 18 July 2023).

Paananen, V. et al. (2023). Using Text-to-Image Generation for Architectural Design Ideation. In *arXiv*. <<https://arxiv.org/abs/2304.10182>> (accessed 18 July 2023).

Ploennings, J., Berger, M. (2022). AI Art in Architecture. In *arXiv*. <<https://arxiv.org/abs/2212.09399>> (accessed 18 July 2023).

Prix, W. et al. (2022). The Legacy Sketch Machine. From Artificial to Architectural Intelligence. In *AD, Machine Hallucinations: Architecture and Artificial Intelligence*, V. 92, No. 3, pp. 14-21.

Ramesh, A. et al. (5 January 2021). DALL-E: Creating images from text. <<https://openai.com/research/dall-e>> (accessed 18 July 2023).

Radford, A. et al. (2021). Learning Transferable Visual Models from Natural Language Supervision. In M. Meila, T. Zhang (eds.). *Proceedings of the 38th International Conference on Machine Learning: Virtuale*, 18-24 July, V. 139, pp. 8748-8763. Maastricht: ML Research Press.

Reed, S. et al. (2016). Generative Adversarial Text to Image Synthesis. In M. F. Balcan, K. O. Weinberger (a cura di). *Proceedings of the 33rd International Conference on Machine Learning*, V. 48, pp. 1060-1069. Maastricht: ML Research Press.

Rombach, R. et al. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, New Orleans, 18-24 June, pp. 10674-10685. New York: IEEE.



Singer, U. et al. (2022). *Make-A-Video: Text-to-Video Generation without Text-Video data*. In *arXiv*. <<https://arxiv.org/abs/2209.14792>> (accessed 18 July 2023).

Tong, H. et al. (2023). An attempt to integrate AI-based techniques into first year design representation course. In K. Vaes, J. Verlinden (a cura di). *Connectivity and Creativity in times of Conflicts. Cumulus Conference Proceedings*. Anversa, 12-15 April, pp. 1-5. Antwerp: University of Antwerp.

Tsigkari, M. et al. (29 March 2021). Towards Artificial Intelligence in Architecture: How machine learning can change the way we approach design. In *Plus Journal*, <<https://www.fosterandpartners.com/insights/plus-journal/towards-artificial-intelligence-in-architecture-how-ma->

[chine-learning-can-change-the-way-we-approach-design](#)> (accessed 18 July 2023).

Wallish, S. (2022). GAN Hadid. In S. Carta (ed.). *Machine Learning and the City: Applications in Architecture and Urban Design*, pp. 477-481. Hoboken-Chichester: John Wiley & Sons.

Yildirim, E. (2022). Text-to-image generation A.I. in architecture. In H. Hale Kozlu (ed.). *Art and Architecture: Theory, Practice and Experience*, pp. 97-119. Lyon: Livre de Lyon.

Zhang, L., Agrawala, M. (2023). Adding Conditional Control to Text-to-Image Diffusion Models. <<https://arxiv.org/abs/2302.05543>> (accessed 18 July 2023).