

Deep Learning and Multimedia Information Analysis
MSc in Artificial Intelligence
Aristotle University of Thessaloniki

Homework 1 & 2

Deep Convolutional Neural Networks & Transformers

**DETR : End-to-End Object Detection with
Transformers**



George Kalitsios
Student Number 62
Thessaloniki 30/5/2021

1. Εισαγωγή

Δεν υλοποιήθηκε εργασία πάνω μόνο σε Βαθιά Συνελικτικά Νευρωνικά Δίκτυα(αν και η αρχική σκέψη ήταν να υλοποιηθεί ένας Faster R-CNN ή Mask R-CNN) αλλά συνδυάστηκαν αυτά μαζί με transformers. Για τις πρώτες δυο εργασίες χρησιμοποιήθηκε ο DETR (Detection Transformer) για λυθεί ένα Object Detection task. Η ιδέα ήρθε στην πρώτη παρουσίαση του μαθήματος όπου έγινε αναφορά στον DETR και διαβάζοντας στην συνέχεια το paper και ψάχνοντας παραπάνω την συγκεκριμένη αρχιτεκτονική μου φάνηκε αρκετά ενδιαφέρον να χρησιμοποιηθούν οι transformers που έχουν κυριαρχήσει στο NLP σε ένα παραδοσιακό Computer Vision task όπως είναι το Object Detection. Εκτός από το paper (DETR : End-to-End Object Detection with Transformers) μελετήθηκαν ακόμα τα Attention Is All You Need και BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding για την καλύτερη κατανόηση των transformers και concept όπως Encoder, Decoder, Self Attention, Multi-head attention αλλά και 2 νέα paper που έρχονται να βελτιώσουν επιπλέον την αρχική υλοποίηση του Facebook AI το (Deformable DETR: Deformable Transformers for End-to-End Object Detection) που βελτιώνει την απόδοση στα small objects και το (Efficient DETR: Improving End-to-End Object Detector with Dense Prior) που βγήκε πριν από έναν μηνά (Απρίλιο του 2021) .

2. Dataset

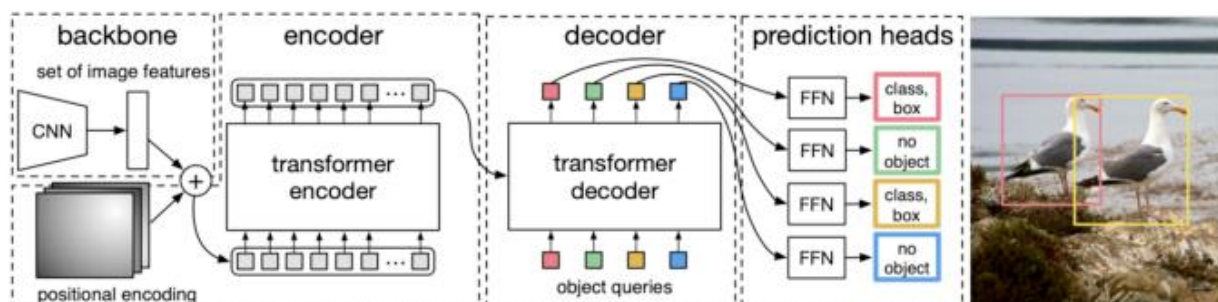
WIDER FACE: A Face Detection Benchmark

Link : <http://shuoyang1213.me/WIDERFACE/>

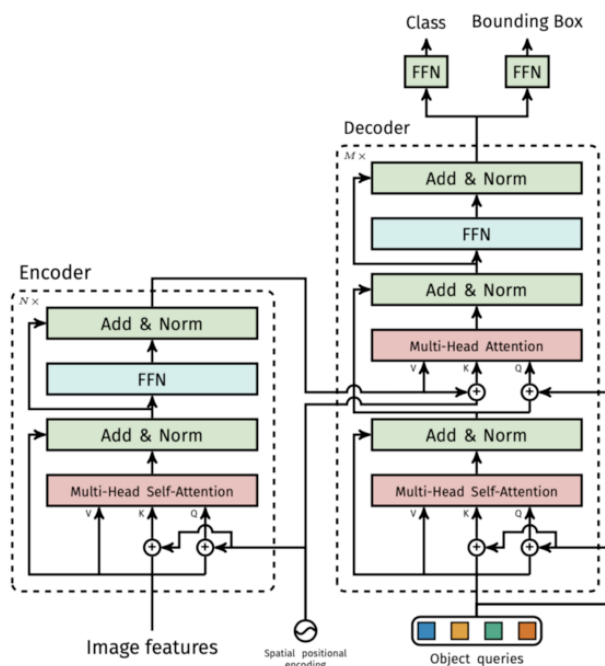
Το dataset έχει 32.203 εικόνες και συνολικά 393.703 πρόσωπα αλλά δεν είναι σε coco format τα annotations. Στον κώδικα που φέρνει το dataset σε coco format κρατήσαμε μόνο εικόνες που έχουν μέχρι 10 πρόσωπα και όχι παραπάνω για να μπορούμε και να ελέγξουμε καλύτερα τον αριθμό των object queries στην αρχιτεκτονική. Επίσης στα πειράματα με τις παραμέτρους δεν χρησιμοποιήθηκε ολόκληρο το dataset αλλά υποσύνολο αυτού ώστε να μπορεί να γίνει πιο γρήγορα ο πειραματισμός καθώς με ολόκληρο το dataset ο χρόνος ήταν απαγορευτικός.

3. Αρχιτεκτονική του DETR

DETR Architecture



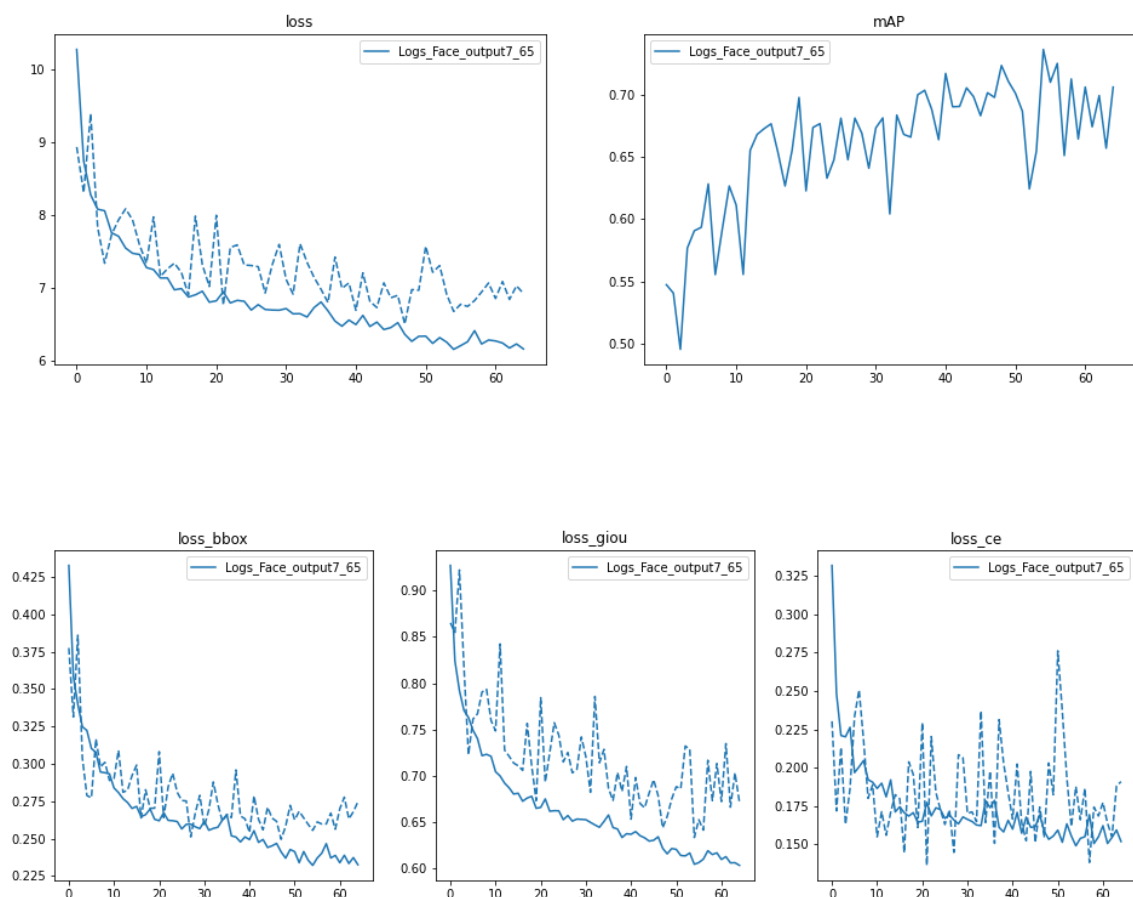
DETR Transformer



4. Πειράματα

- Τα πειράματα εκτελέστηκαν τοπικά σε μια Nvidia RTX 2080 Ti
- Αρχικά δοκιμάστηκε με λίγα δεδομένα να δούμε αν το μοντέλο κάνει overfitting καθώς αναφέρθηκε και στο μάθημα ότι πρέπει να είναι από τα πρώτα που μπορούμε να δοκιμάσουμε όταν κάνουμε debugging.
- **Πείραμα 1ο - Using the entire Dataset:**

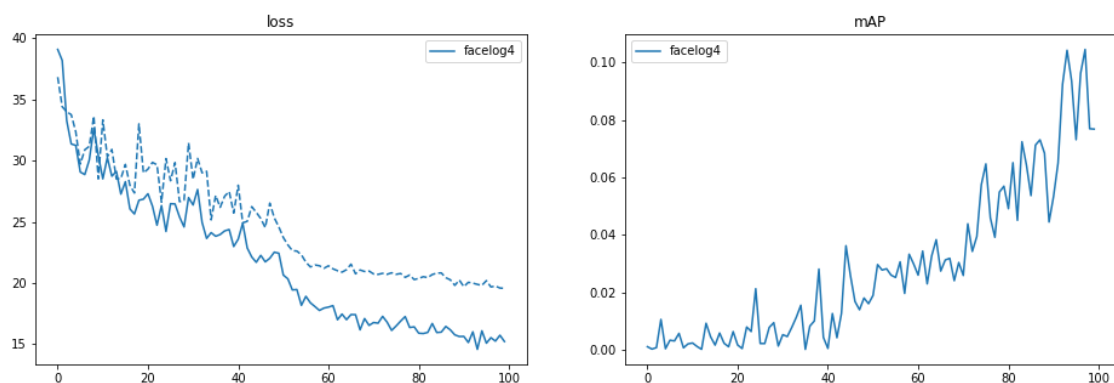
Στο πείραμα αυτό Χρησιμοποιήθηκε ολοκληρωτο το Dataset ο χρόνος εκπαίδευσης για 65 Epochs ήταν 32,5 ώρες καθώς για να δεις τις 10000 εικόνες του training set χρειαζόταν περίπου 30 λεπτά για 1 epoch. Το Map (Mean Average Precision) αυξάνεται καθώς προχωράει το training και το loss μειώνεται, παρακάτω φαίνονται επίσης και τα 3 διαφορετικά losses [classification loss, L1 bounding box distance loss, and GloU (Generalized Intersection over Union) loss].



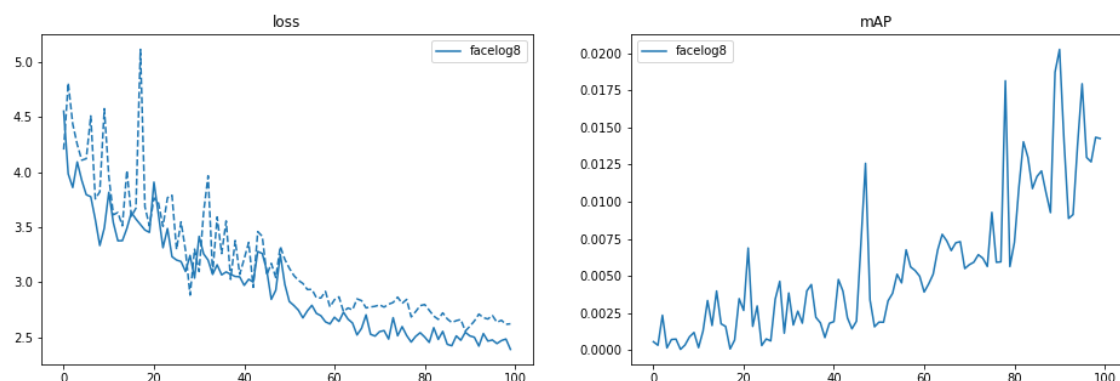
Επειδή το Dataset ήταν αρκετά μεγάλο και ο χρόνος εκπαίδευσης ήταν απαγορευτικός στα πειράματα με τις διαφορετικές παραμετρους επιλεχθηκε ένα υποσυνολο με 400 εικονες.

- **Πείραμα 2ο - Number of encoder/decoder layers:**

Ας ξεκινήσουμε βλέποντας κάποια πειράματα που δεν ήταν καλά και τα συμπεράσματα που προκύπτουν από αυτά. Δοκιμαστηκε να αλλάχθει στην αρχιτεκτονική ο αριθμός των encoders και ο αριθμός των decoders. Παρακατω βλέπουμε το loss και το mAP αυξανοντας από 6 σε 8 των αριθμο για τους encoders και τους decoders και βλέπουμε ότι το Map Μετα από 100 Epochs είναι πολύ κακο(10%). Στο επίσημο paper υπηρχαν πειράματα που με την αύξηση των encoders/decoders βελτιωναν τα αποτελεσματα εδώ αυτό δεν συμβαίνει γιατί έχουμε ένα υποσυνολο από το dataset ενώ στο paper χρησιμοποιούσαν την Coco.

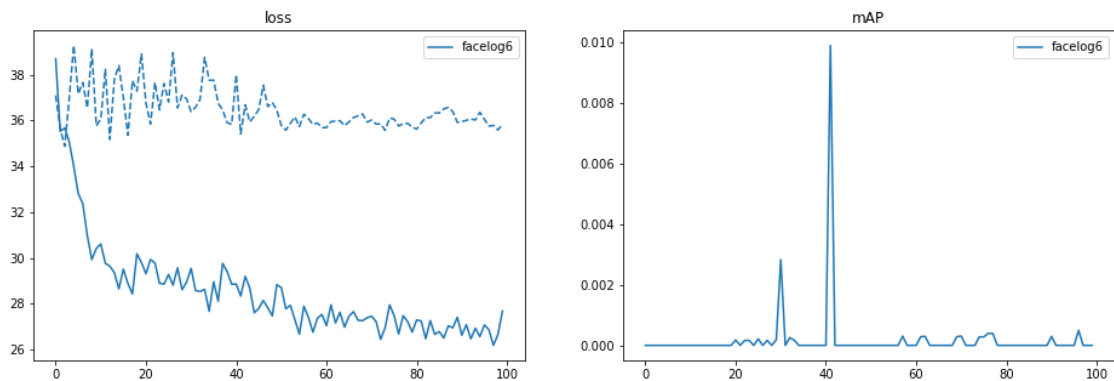


Παρακατω βλέπουμε για 3 encoders layers και 1 decoder και βλέπουμε ότι ενώ το loss είναι χαμηλο σε σχέση με πριν το mAP παλι δεν είναι καθολου καλο σχεδον 0.



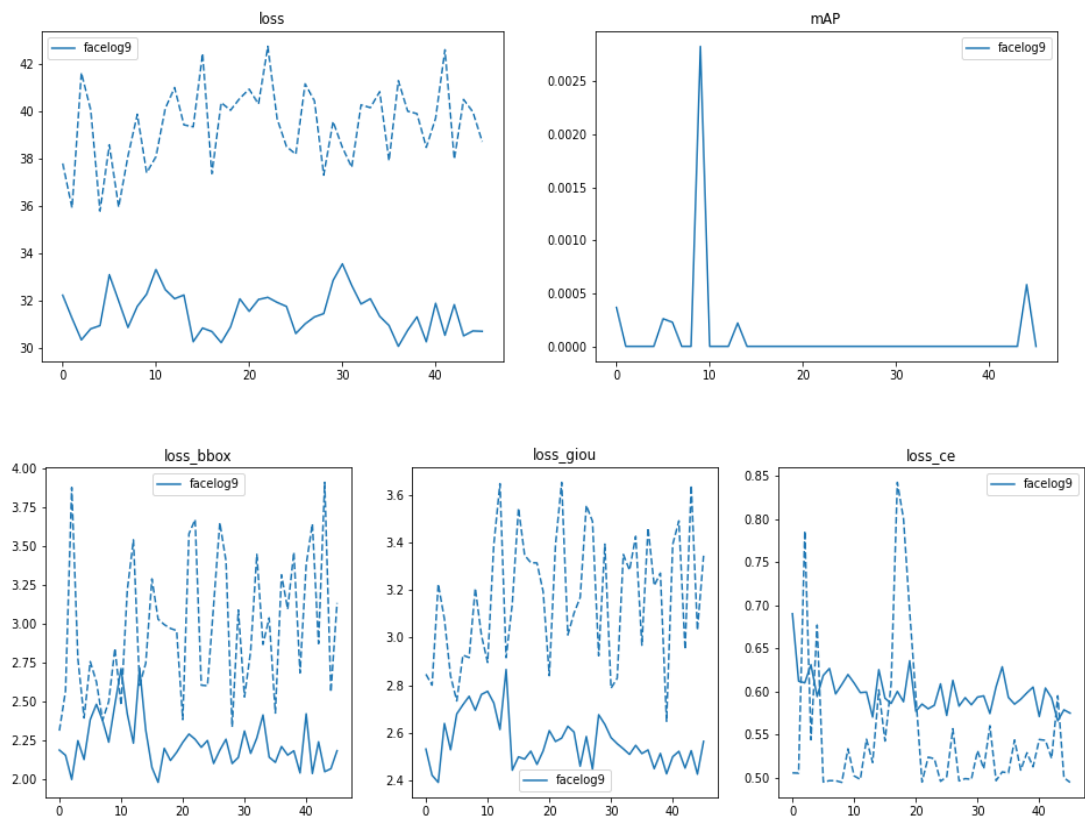
- **Πείραμα 3ο - Training from scratch:**

Το να εκπαιδευσουμε from scratch με 400 εικόνες δεν ήταν καλή ιδέα και είναι και λογικό όπως φαίνεται παρακατω. Αν χρησιμοποιούσαμε ολοκληρο το dataset θα μπορούσαμε να το κανουμε αφού θα ειχαμε τουλαχιστον 10000 εικονες, για αυτό χρησιμοποιησαμε pre-trained βαρη από την Coco και καναμε fine-tuning.



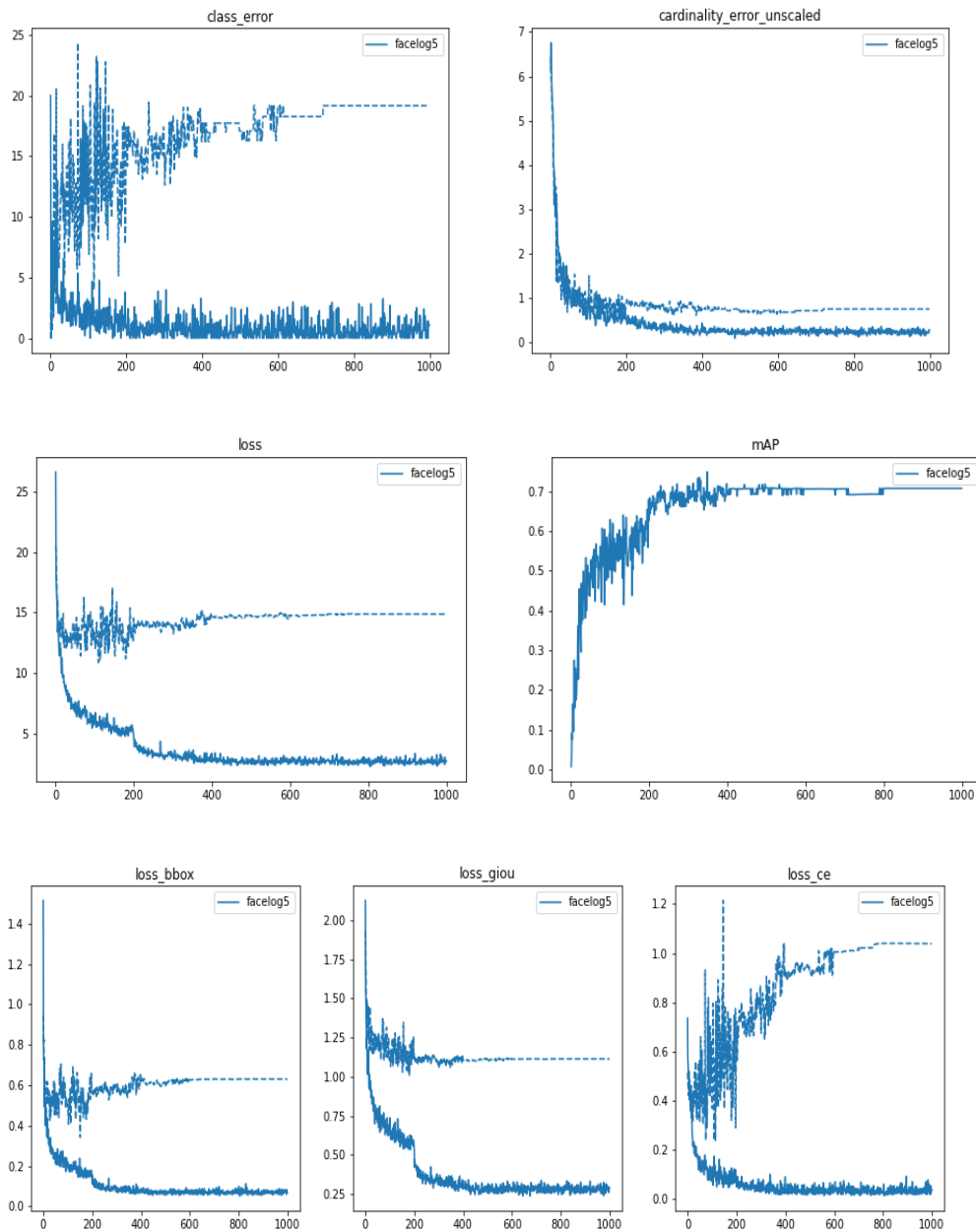
- **Πείραμα 4ο - μεγάλο learning rate=10-3:**

Εδώ βλέπουμε άλλο ένα κακο πειραμα που χρησιμοποιηθηκε μεγαλο learning rate και από ότι φαίνεται αυτό δεν βοήθησε καθολου το δικτυο,πιεζοντας δηλαδή να μαθει γρηγορα δεν εχει τα επιθυμητα αποτελεσματα.



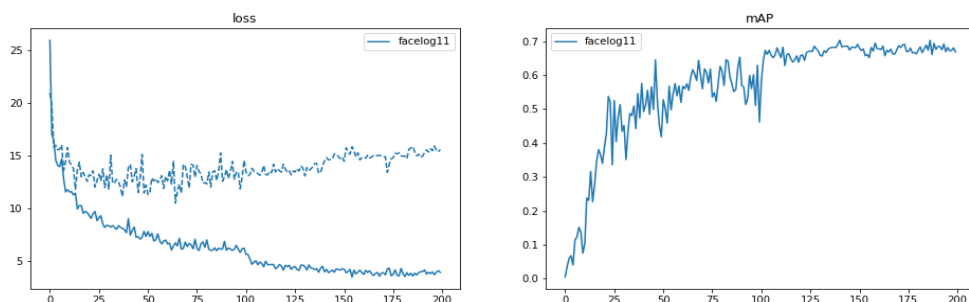
- **Πείραμα 5ο – Training 1000epochs:**

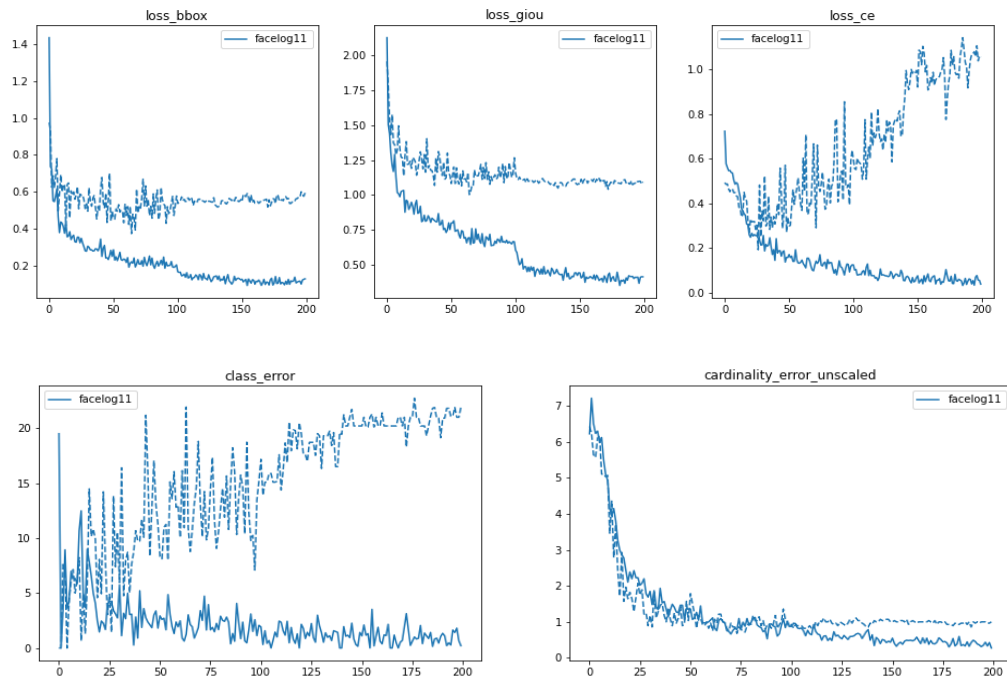
Βλεπουμε εδώ ένα πειραμα για 1000 epochs που το Overfitting γινεται πολύ Εντονο μετα τα 150 epochs.



- **Πείραμα 6ο – Training 200epochs:**

Αναλογη συμπεριφορα βλεπουμε και για 200 epochs σε σχεση με πριν

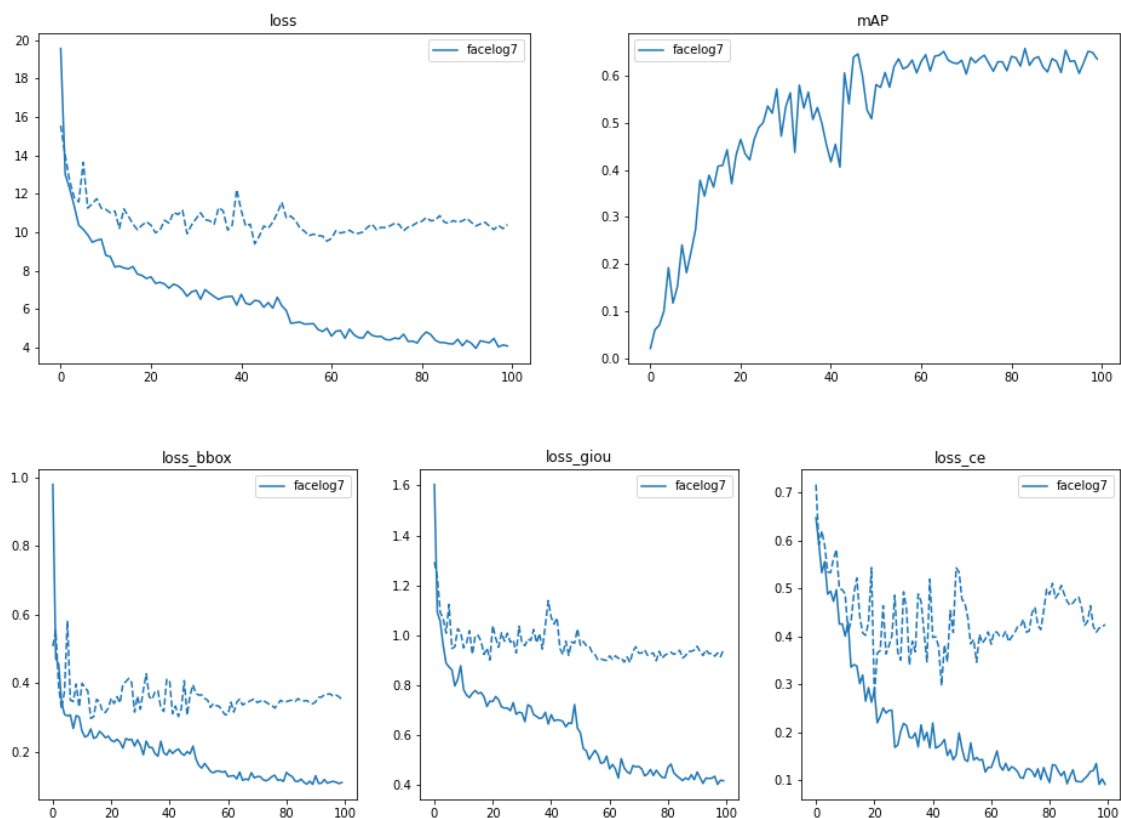




Θα μπορούσαμε να το αντιμετωπίσουμε εκπαιδευοντας με παραπάνω δεδομένα.

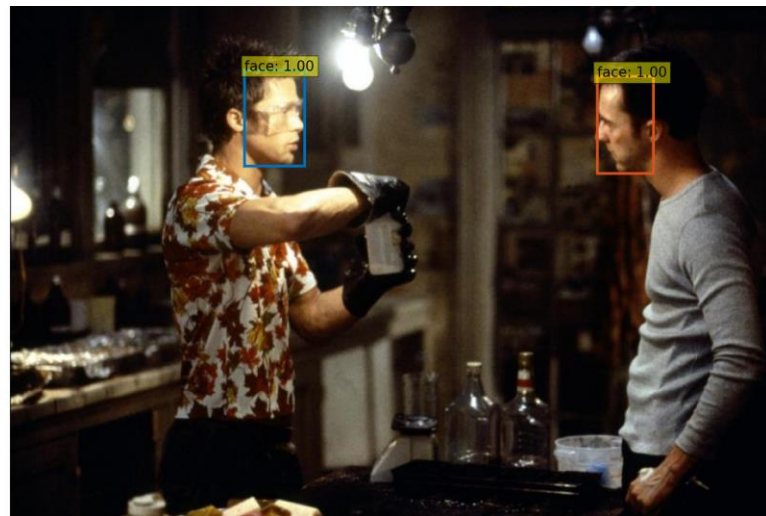
- **Πείραμα 7ο – Training Number of queries=50:**

Παρατηρούμε εδώ ότι αυξανοντας τα queries από 10 σε 50 το loss πηγαινει καλυτερα.



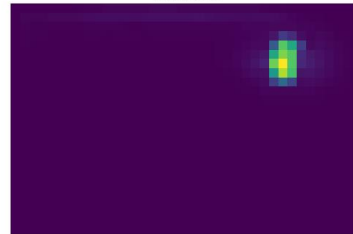
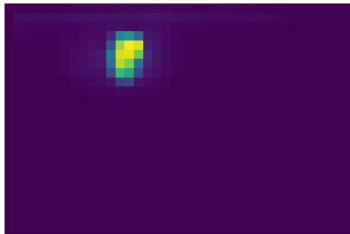
5. Visualizing predictions - Visualize encoder-decoder multi-head attention weights

•



query id: 1

query id: 5



•



query id: 0

query id: 1

query id: 2

query id: 3

query id: 4

query id: 5

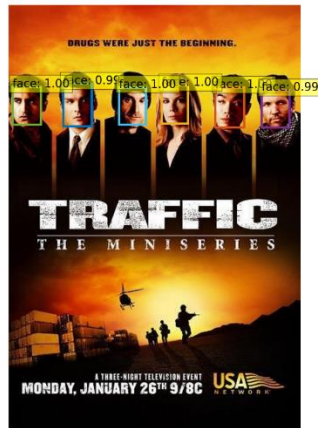
query id: 6

query id: 7

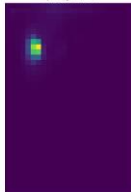
query id: 8

query id: 9





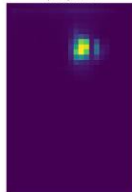
query id: 0



query id: 2



query id: 5



query id: 7



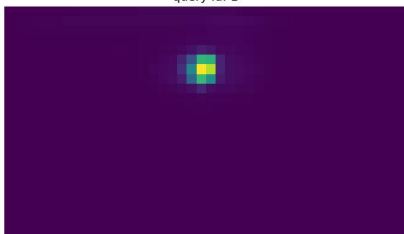
query id: 8



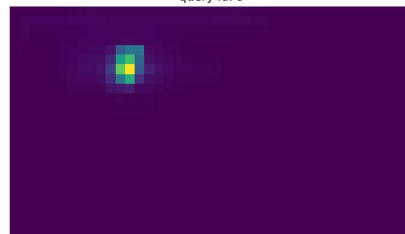
query id: 9

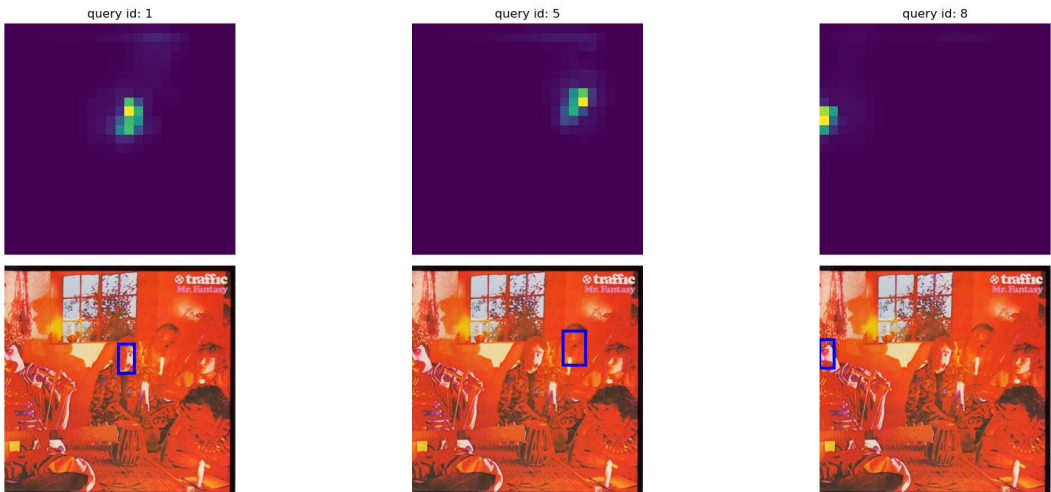
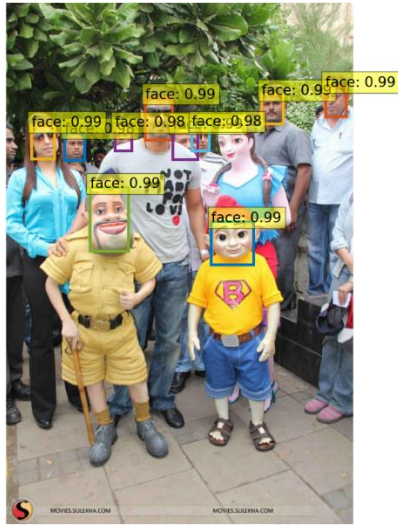


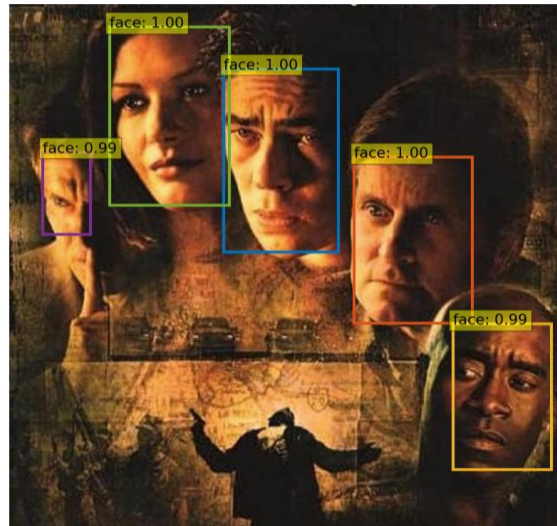
query id: 1



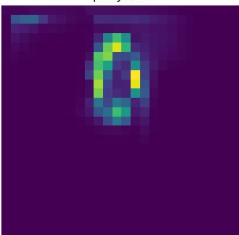
query id: 9



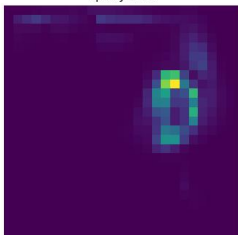




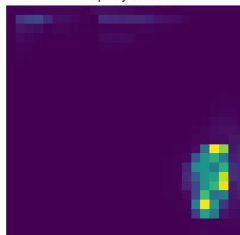
query id: 1



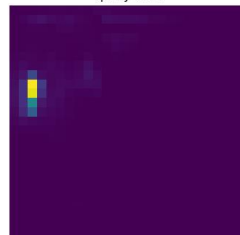
query id: 5



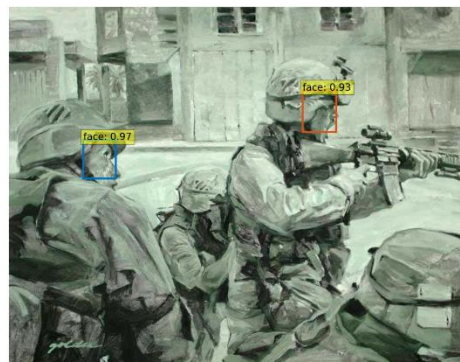
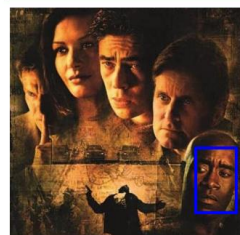
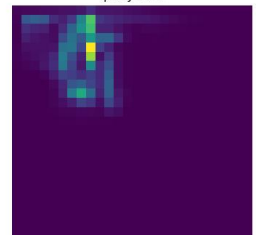
query id: 6



query id: 8



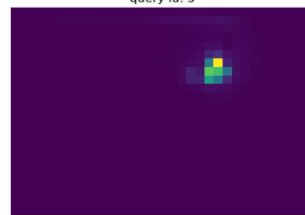
query id: 9



query id: 0



query id: 5



6. Επίλογος

Στην εργασία αυτήν χρησιμοποιήθηκε μια state-of-the-art Αρχιτεκτονική η οποία προτάθηκε μέσα στο 2020 και εισάγει transformers σε ένα κλασσικό πρόβλημα του Computer Vision όπως είναι το Object Detection, σημαντικό πλεονέκτημα της είναι η απλότητα της και ότι δεν υπάρχουν πλέον τα Anchors & Non-maximum suppression components που υπάρχουν σε κλασσικές αρχιτεκτονικές για object detection (Mask R-CNN, Faster κλπ.) και πρέπει να γίνει σωστό tuning των παραμέτρων τους. Λύσαμε ένα Face Detection πρόβλημα, χρησιμοποιήθηκε ένα custom Dataset με πρόσωπα και έγινε fine-tuning ξεκινώντας με pretrained βάρη από την COCO. Το dataset έπρεπε να έρθει σε coco format, να αλλάξει ο data loader και να γίνουν και αλλαγές σε άλλα σημεία στον αρχικό κώδικα ώστε να μπορεί να χρησιμοποιηθεί για το πρόβλημα του face detection. Σημαντικός χρόνος αφιερώθηκε στην μελέτη του αρχικού paper (DETR: End-to-End Object Detection with Transformers) αλλά και στην κατανόηση των Transformers. Επιπλέον έγιναν αρκετά πειράματα αλλάζοντας παραμέτρους όπως είναι το learning rate, learning rate drop, τον αριθμό των encoders και decoders της αρχιτεκτονικής, τον αριθμό των number queries κλπ. και προσπαθώντας κάθε φορά να ερμηνευτεί το αποτέλεσμα των αλλαγών μας. Παρουσιάστηκαν ακόμα οπτικά αποτελέσματα με της προβλέψεις του μοντέλου αλλά και τα attention weights του τελευταίου layer του decoder και είδαμε ποια pixel στην εικόνα ήταν αυτά που οδήγησαν το μοντέλο να κάνει μια συγκεκριμένη πρόβλεψη.