

A. MODEL SUMMARY

A1. Background on you

Competition Name: RSNA 2024 Lumbar Spine Degenerative Classification
(<https://www.kaggle.com/competitions/rsna-2024-lumbar-spine-degenerative-classification>)

Team Name: HLIP

Private Leaderboard Score: 0.401322

Private Leaderboard Place: rank 7th

A2. Background on you/your team

What is your academic/professional background?

I am a contract computer vision and deep learning algorithm engineer. My job includes discovering fracture in x-ray images and implementing visual slam for robotic navigation. Recently, I am helping companies to fine-tune LLM models for their applications.

Did you have any prior experience that helped you succeed in this competition?

I am familiar with deep learning and build deep models in my work. I took part in previous RSNA competitions before like RSNA 2023 Abdominal Trauma Detection and RSNA 2022 Cervical Spine Fracture Detection. These are helpful for me.

What made you decide to enter this competition?

The problem of analyzing 3d MRI volume scan is interesting for me. I am generally interested in medical image analysis and computer vision problems.

How much time did you spend on the competition?

About 5 hours per day, over two months.

A3. Summary

In this competition, our task is to predict lumbar spine degenerative conditions from given 3d volume MRI scans. Since my teammate @lihaoweicvch already has a very good two-stage model, my task is to design a one-stage model to improve his results.

While a two-stage model uses "crop and classify" method, a one-stage model uses "point-masking and pool" instead. Figure.1 and Figure.2 show deep net architecture for NFN (neural foraminal narrowing) and SCS (spinal canal stenosis) detection respectively. In summary, the main processing layers for both architectures are:

1. Predict the 5 lumbar level point pixel-wise heatmap for the L1/L2 to L5/S1. For learning the net parameter, we use the Jensen-Shannon divergence loss function.
2. Convert point pixel-wise heatmap to target xyz coordinates using the differentiable spatial-to-numerical transform (DSNT) from in paper [1].
3. Predict the 3 spine degeneration grade pixel-wise probability for the classes Normal/Mild, Moderate, or Severe.

4. Use the point pixel-wise heatmap to pool the grade pixel-wise probability using multiplication and summing over volume or area. The pooled results is the target level-wise grade conditions for submission,

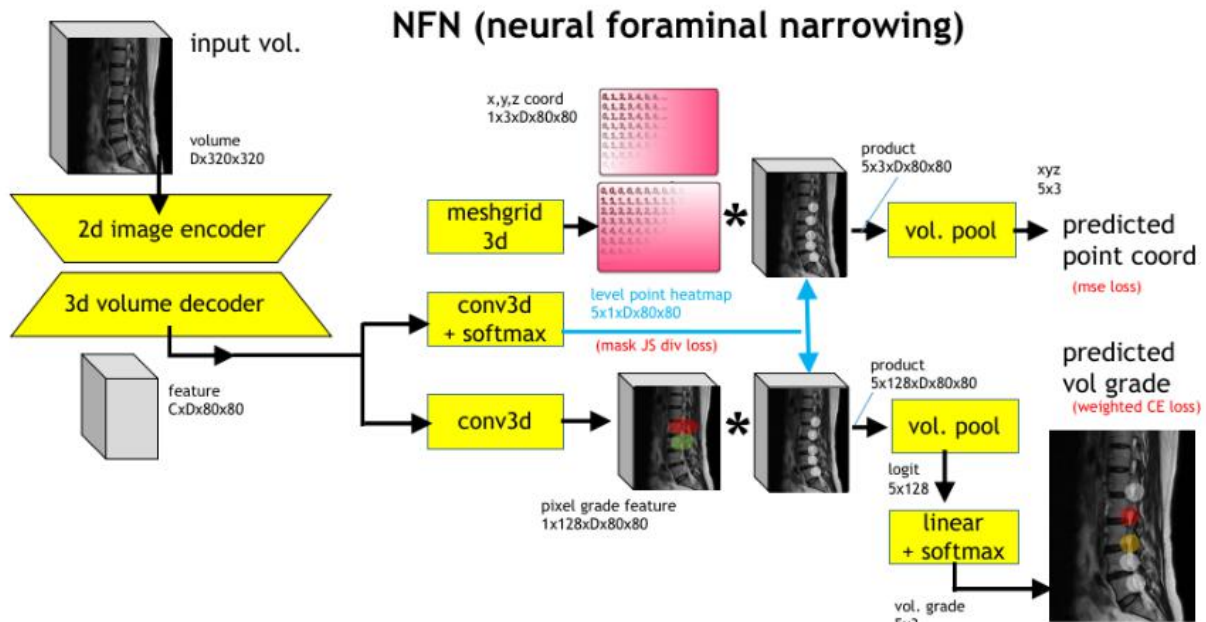


Figure.1 Deep net architecture for NFN (neural foraminal narrowing) detection

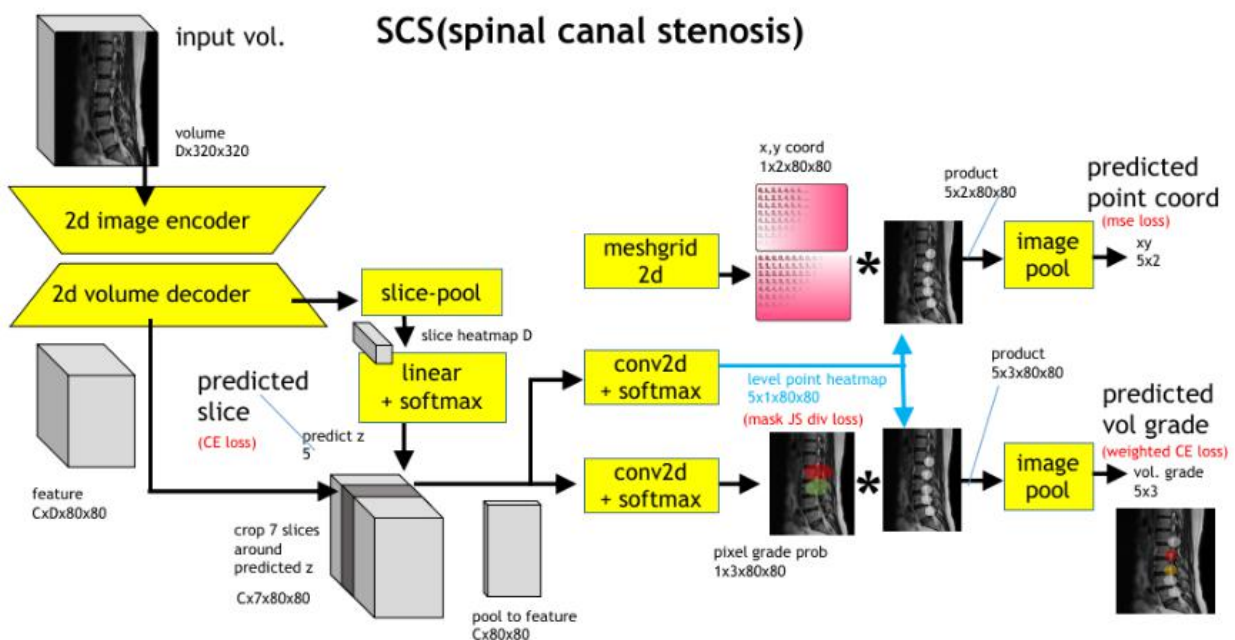


Figure.2 Deep net architecture for SCS (spinal canal stenosis) detection

A4. Features Selection / Engineering

1. 1. Image encoder

We are using pyramid transformer PVTv2 [2] as image encoder backbone for submission. They have very large context which we think are important to get good results. For example, both left and right level are likely to have the same spine degeneration conditions.

2. Point-masking and pooling

This is first mentioned in Section A.3 Summary. The level point heatmap can be seen as learned attention map to select with pixel we should look at for spine degeneration conditions.

A5. Training Method(s)

Training is straight forward with multi-task loss setup for each net in Figure.1 and 2:

- pixel-wise level point : Jensen-Shannon divergence loss
- pixel-wise grade probability : Cross entropy loss
- volume-wise level grade probability : weighted cross entropy loss

We perform back propagation with ADAM optimizer. The training hyper-parameters varies for different net due to the difficulty of the problem. Learning rate is from $1e-4$ to $1e-5$ and training epochs ranges from 50 to 100. Batch sizes is from 3 to 8 volumes.

For hardware, we use one Nvidia Ada A6000/48GB Ampere GPU. Powerful GPU makes our training fast and large memory enables us to train with whole 3d volume.

Some observations from training are:

- We can temporarily modify the network to use the ground truth level point heatmap and focus solely on learning grade prediction. This serves as an upper bound for grade prediction performance.
- One-stage networks are challenging to train because the optimal parameters for coordinate prediction are not necessarily the best for grade prediction. This can be confirmed by observing the loss for level and grade during training iterations. They seem to be competing with each other.

We believe one reason for this is labeling errors, particularly the significant confusion between L5 and S1 level labels. Additionally, the point labels themselves are not very precise.

We had a bug in flip augmentation at train. We forgot to reorder the left right level points (i.e. at training the left is mistakenly labeled as right and vice versus for 50% of the time). We discovered the bug just one day before submission and was unable to correct all models for submission. Hence the final submission consisted of both bugged and unbugged models.

We later make post submission after the competition deadline with unbugged models only and showed better results. Next, we summarize our observation:

1. One-stage NFN model indeed improves results. The baseline two stage has private /public score of 0.406/0.3500.
2. Even with bugged models only, the improvement is 0.402/0.3481.
3. For unbugged models, we have the best score of 0.400/0.3466.
4. One stage SCS models did not improve results and are not included for final submission. The reasons are unclear. When alone, one stage SCS model are have comparable accuracy as the two-stage counterparts. However after ensemble, the combined results falls a little.

What was the most important trick you used?

One stage model improves two stage model in ensemble. Blending for submission is:
 $0.75 \times 2\text{stage} + 0.25 \times 1\text{stage}$

What do you think set you apart from others in the competition?

None of the top winning team uses one-stage model.

A7. Simple Features and Methods

Many customers are happy to trade off model performance for simplicity. With this in mind:
Is there a subset of features that would get 90-95% of your final performance?

From our report it is sufficient just to use only two stage mode. There will be drop of 5% accuracy though.

A8. Model Execution Time

How long does it take to train your model?

Training takes 6 hours for one net for each one fold, using Nvidia Ada A6000 Ampere GPU with memory 48 GB.

How long does it take to generate predictions using your model?

For one stage model, It take less than a minute to process one volume of MRI scan.

How long does it take to train the simplified model (referenced in section A7)?

I expect a reduction of more than 50% in time.

How long does it take to generate predictions from the simplified model?

I expect a reduction of more than 50% in time.

A9. References

[1] "Numerical Coordinate Regression with Convolutional Neural Networks" - Aiden Nibali, Arvix 2018, <https://arxiv.org/abs/1801.07372>

[2] "PVT v2: Improved Baselines with Pyramid Vision Transformer" - Wenhai Wang
<https://arxiv.org/abs/2106.13797>