

# Aplicação de Patches Dinâmicos em Vision Transformers em Exames de Papanicolau

Vinícius H. Giovanini  
Alexei Machado

Pontifícia Universidade Católica de Minas Gerais

13 de dezembro de 2024

# Sumário

- 1 Introdução
- 2 Vision Transformer
- 3 Materiais
- 4 Métodos
- 5 Experimentos
- 6 Resultados
- 7 Conclusão
- 8 Referências

# Análise de Imagens Médicas com DeepLearning

- Classificação de imagens médicas (radiografias, tomografias, ressonâncias magnéticas, etc.);
- Redes Neurais Convolucionais (CNNs);
- Visions Transformers (ViTs);
- Modelos Híbridos;
- Análise automática em grande escala;
- É essencial alcançar um elevado índice de acerto.

# O Exame de Papanicolau

- Exame de Papanicolau surgiu em 1928 com George Papanicolau;
- Detectar precocemente o câncer de colo de útero;
- Redução significativa na mortalidade por câncer cervical;
- Incidência de 4,86 casos por 100.000 mulheres no Brasil [Barcelos et al., 2017].

# Fluxo de Entrada e Saída do ViT

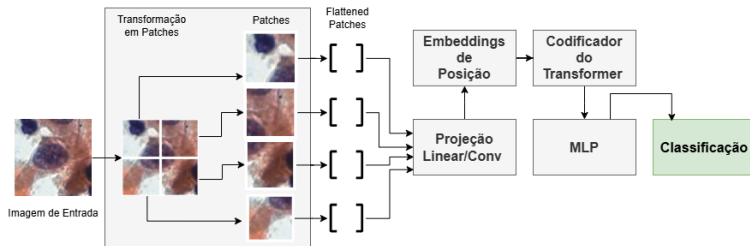


Figura: Fluxo de entrada do ViT

# Objetivos: ViT com Patches Dinâmicos

- Utilização do Vision Transformer;
- Base de Dados Centro de Reconhecimento e Inspeção de Células (CRIC) - UFOP
- Realização do Fine Tuning;
- Incorporação dos Patches Dinâmicos ao ViT;
  - ① Seleção Randômica (SR);
  - ② Seleção Randômica Aprimorada (RA);
  - ③ Seleção por Segmentação (SS);
- Comparação entre ViT e CNN com ensemble [N. Diniz et al., 2021].

- Improving robustness for vision transformer with a simple dynamic scanning augmentation [Kotyan and Vargas, 2024]
- An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale [Dosovitskiy et al., 2021]
- Visual Transformers: Token-based Image Representation and Processing for Computer Vision [Wu et al., 2021]
- Cric searchable image database as a public platform for conventional pap smear cytology data [Rezende et al., 2021]
- A Deep Learning Ensemble Method to Assist Cytopathologists in Pap Test Image Classification [N. Diniz et al., 2021]

# Arquitetura do Vision Transformer (ViT)

- Introduzido em 2021 pelo Google Research;
- Baseado no Transformers (NLP);
- Arquitetura é composta:
  - 1 Divisão de imagens em patches;
  - 2 Adição do CLS Token e do Embeddings de Posição;
  - 3 Processamento pelo códecificador (MLP e Auto-Atenção);
  - 4 MLP de Classificação.



# Arquitetura do Vision Transformer (ViT)

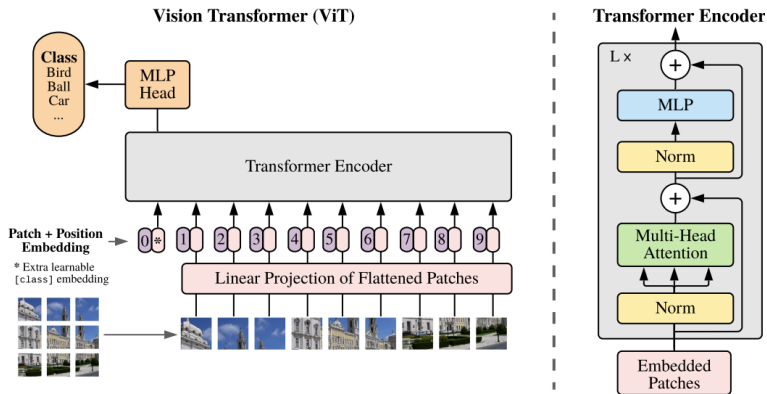


Figura: Arquitetura do ViT [Dosovitskiy et al., 2021]

# Patches

- Divisão de uma imagem em blocos menores de tamanho fixo;
- Cada patch é tratado como um token;
- Normalmente utiliza-se tamanho de 16x16 e 32x32;
- Os patches são mapeados para vetores;
- Menos recursos computacionais.

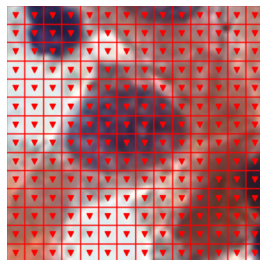


Figura: ASC-H 10

$$\text{Total de patches} = \left( \frac{H}{PS} \right) \times \left( \frac{W}{PS} \right)$$

## • **Projeção Linear**

- Aplica-se uma camada Linear;
- A imagem é dividida em pequenos patches, que são então achatados em vetores unidimensionais;
- Utilizada no trabalho [Dosovitskiy et al., 2021];
- Própria imagem.

## • **Projeção Convolutacional**

- Filtro convolutacional para gerar os embeddings;
- Utilizada no trabalho [Wu et al., 2021];
- Preserva-se informações espacial local;
- Mapa de características.

# Materiais - Conjunto de Dados (Dataset)

- Centro de Reconhecimento e Inspeção de Células (CRIC);
- 400 imagens de exames de Papanicolau;
- Arquivo CSV e JSON;
- 11.534 núcleos de célula;
- Sistema Bethesda.
  - ① Atypical squamous cells of undetermined significance (ASC-US)
  - ② Atypical squamous cells cannot exclude a high-grade lesion (ASC-H)
  - ③ High-grade squamous intraepithelial lesion (HSIL)
  - ④ Low-grade squamous intraepithelial lesion (LSIL)
  - ⑤ Negative for intraepithelial lesion (NFIL)
  - ⑥ Squamous cell carcinoma (SCC)

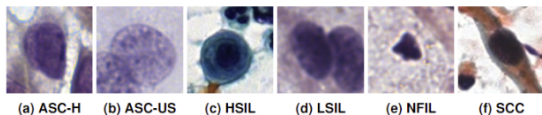


Figura: Recorte das células 90x90

# Balanceamento de dados

- Grande desbalanceamento entre os dados;
- 80/20 e 20% do Treino para o Teste;
- Albumentations;
- Corte Aleatório, Inversão Horizontal e Rotação;
- Balanceamento para 1000 imagens no Treino.

Divisão	ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
Treino (A/D)	592/1000	388/1000	1.090/1000	871/1000	4.339/1000	103/1000
Validação	185	122	341	272	1.356	33
Teste	148	96	272	217	1.084	25

**Tabela:** Quantidade de Imagens antes e depois do balanceamento

# Método - Introdução aos Patches Dinâmicos

- Explorar a sobreposição de patches;
- Extrair área de interesse;
- Substituição da classe tradicional de Embeddings por uma classe customizada;
- Introdução ao conceito de centros.

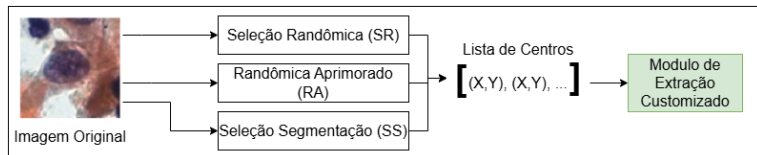


Figura: Fluxo Inicial de Extração de Patches com Centros

# Método - Seleção Randômica (SR)

- Método de extração completamente aleatório;
- Seleciona X centros a partir da quantidade total de patches necessários;
- Método rápido, possível utiliza-lo durante o treinamento;
- Sobreposição descontrolada.

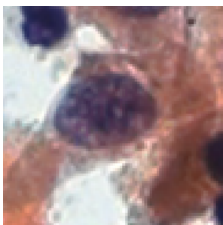
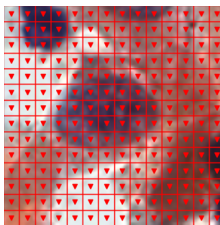
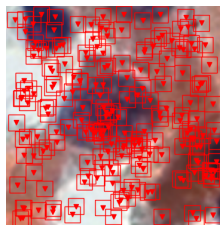


Imagem Original



Extração por Grid



Extração por SR

# Método - Randômico Aprimorado

- Método de extração aleatório;
- Algoritmo que controla sobreposição;
- Centro só é validado se não estiver dentro da área de outro patch já existente;
- Método custoso, não é realizado durante o treinamento;
- Sobreposição controlada.

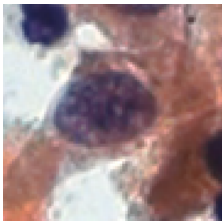
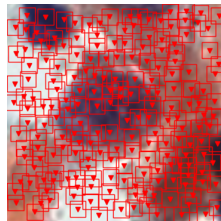


Imagem Original



Extração por SR



Extração por RA



# Método - Seleção por Segmentação

- Método de extração por área de interesse;
- Utiliza-se o GrabCut;
- Extrai os patches em Grid da máscara;
- Patches que sobram são selecionados randomicamente;
- Método custoso, não é realizado durante o treinamento;
- Sobreposição pouco controlada.



Imagem Original

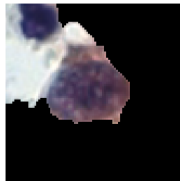


Imagem Recortada



Máscara



SS

# Método - Seleção por Segmentação



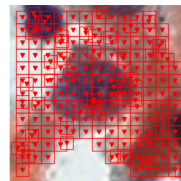
Imagem Original



Imagem Recortada



Máscara



SS

# Método - Pré-processamento dos Patches Dinâmicos

- Pré-processamento para cada abordagem;
- A lista de centros foi salva em um dicionário;
- Classe personalizada de carregamento de dados (ImageFolder);
- Imagem, Rótulo, Nome da Imagem;
- Busca em um dicionário;
- Tempo do treinamento caiu para 3 horas em todos as extrações dinâmicas.

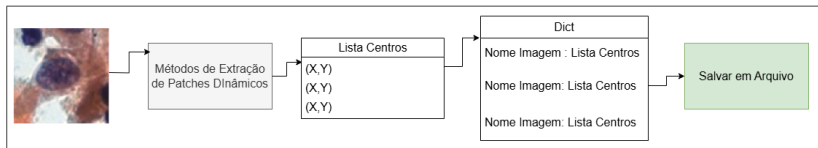
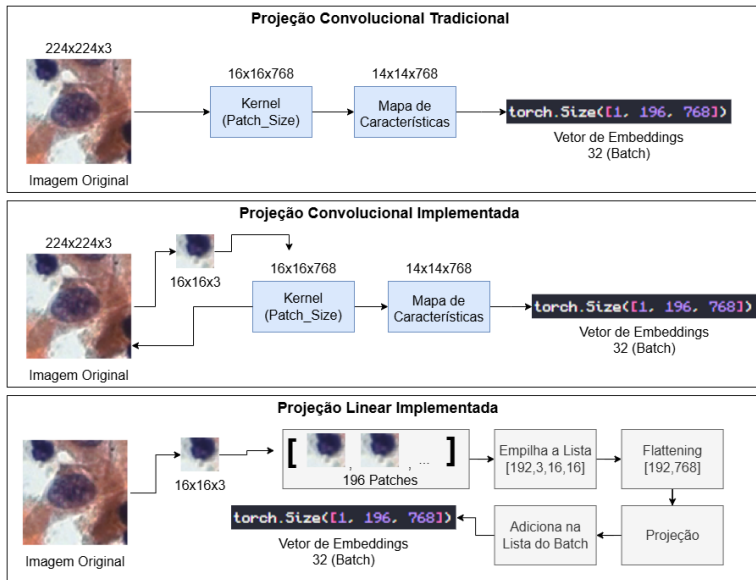


Figura: Fluxo de Pré-Processamento dos Centros

# Método - Modelos pré-treinados de ViT

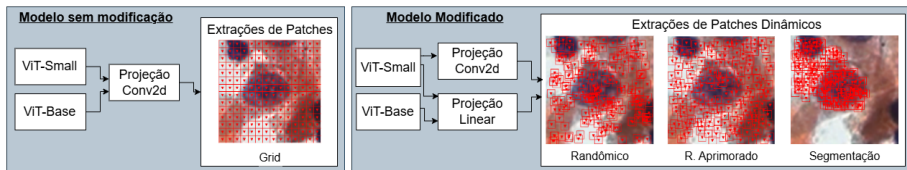
- Google, Facebook (Meta), Microsoft;
- Foi utilizado os modelos do Google;
  - ViT-B/16;
  - ViT-S/16;
  - ViT-T/16;
- Disponibilizados no GitHub do Google Research JAX\FLEX;
- Hugging Face em Pytorch.

# Método - Implementação das Projeções



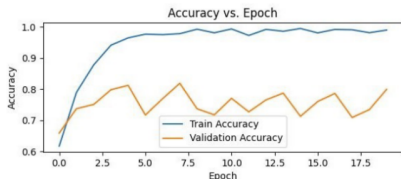
# Fluxo de Experimentos

- 2 treinamentos do modelo com extração por Grid;
  - Projeção Convolutacional;
  - Extração por Grid;
- 9 treinamentos do modelo com extração Dinâmica;
  - Projeção Convolutacional;
  - Projeção Linear.

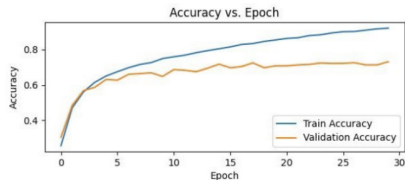


# Experimentos - Hiperparâmetros

- Taxa de Aprendizizado ( $1e-5$ );
- Dropout (0.4);
- Weight decay (Descartado);
- Batch size (32);
- Épocas (60);
- Modelo ViT-Tiny foi descartado.



(a) Taxa de Aprendizizado  $1e-4$



(b) Taxa de Aprendizizado  $1e-5$

Figura: Alteração na Taxa de Aprendizado

# Experimentos - Finetuning - Extração por Grid

- Modelo Base tem 87.6 milhões de parâmetros;
- Modelo Small tem 22.1 milhões de parâmetros;
- Descongelou-se somente o último MLP do último bloco do modelo ViT-Base;
- Descongelou-se 5 blocos completos no modelo ViT-Small.

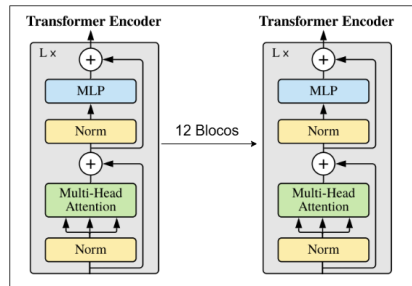
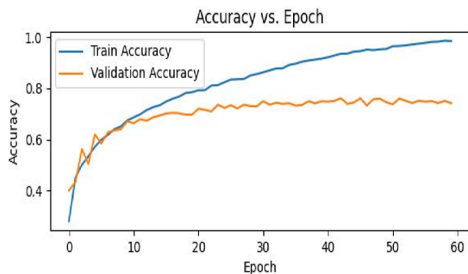


Figura: Bloco do Encoder - ViT



# Gráficos ViT-Base

- 76% (Validação), 75% (Teste)

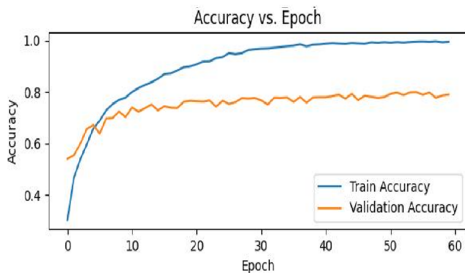


		Confusion Matrix					
		ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
True	ASC-H	74.05	2.16	15.68	3.78	4.32	0.00
	ASC-US	5.74	47.54	4.10	23.77	18.85	0.00
	HSIL	4.40	0.88	84.75	2.05	1.76	6.16
	LSIL	5.88	17.28	4.41	59.93	11.40	1.10
	NFIL	3.83	6.34	2.14	6.56	80.38	0.74
	SCC	6.06	0.00	27.27	3.03	0.00	63.64
		ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
		Predicted					

Figura: Modelo ViT-Base

# Gráficos ViT-Small

- 79% (Validação), 81% (Teste)



		Confusion Matrix					
		ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
True	ASC-H	63.24	4.32	11.89	10.27	9.19	1.08
	ASC-US	2.46	68.85	0.00	12.30	16.39	0.00
	HSIL	2.35	1.76	86.80	3.52	1.76	3.81
	LSIL	1.47	26.47	2.94	55.15	12.87	1.10
	NFIL	0.96	7.74	0.66	3.91	86.65	0.07
	SCC	6.06	0.00	18.18	0.00	0.00	75.76
		ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
		Predicted					

Figura: Modelo ViT-Small

# Experimentos - FineTuning - Extração Dinâmica

- É repetindo as mesmas arquiteturas e hiperparâmetros, para o ViT-Small e ViT-Base;
- Abordagem de extração de patches dinâmica;
- Tipo de projeção Linear e Convolutacional;
- Comparando a acurácia do modelo de extração por Grid.

Confusion Matrix

True \ Predicted	ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
ASC-H	67.03	2.70	16.22	7.03	2.70	4.32
ASC-US	6.56	48.36	0.82	18.85	22.95	2.46
HSIL	11.73	0.29	72.43	2.05	2.93	10.56
LSIL	9.93	12.50	3.31	61.76	9.19	3.31
NFIL	4.28	8.70	2.36	14.01	69.47	1.18
SCC	6.06	0.00	12.12	0.00	0.00	81.82

(a) SR Base Linear

Confusion Matrix

True \ Predicted	ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
ASC-H	54.59	1.62	20.54	4.86	14.05	4.32
ASC-US	9.02	16.39	4.10	12.30	56.56	1.64
HSIL	12.61	0.29	66.28	2.64	10.26	7.92
LSIL	7.72	7.72	9.93	25.37	47.43	1.84
NFIL	2.88	1.62	4.57	4.28	85.69	0.96
SCC	3.03	0.00	21.21	0.00	3.03	72.73

(b) SS Small Linear

Confusion Matrix

True \ Predicted	ASC-H	ASC-US	HSIL	LSIL	NFIL	SCC
ASC-H	64.32	3.24	15.68	10.81	5.95	0.00
ASC-US	3.28	61.48	2.46	13.93	18.85	0.00
HSIL	3.52	1.17	85.92	5.28	2.35	1.76
LSIL	0.74	16.18	3.68	67.65	11.03	0.74
NFIL	1.25	8.63	1.55	8.55	79.87	0.15
SCC	3.03	0.00	42.42	9.09	0.00	45.45

(c) RA Small Conv

Figura: Matriz de Confusão - Melhores Resultados Dinâmicos

- Método CNN com Ensemble obteve acurácia de 95%.

ViT	Projeção	Seleção Randômica	Randômica Aprimorada	Seleção por Segmentação
Base	Linear	69.54	64.92	59.50
Small	Linear	63.73	70.14	70.73
Small	Conv2d	75.02	77.30	75.78

**Tabela:** Acurácia dos Resultados Finais - Extração de Patches Dinâmicos

# Análise dos Resultados

- Nenhuma abordagem de extração, superou a CNN com ensemble;
- A extração de patches dinâmicos impactou a acurácia;
- A projeção Convolutacional apresenta vantagens em relação a Linear;
- O melhor método dinâmico foi o Randômico Aprimorado com Projeção Convolutacional (77.3% acc);
- Esse conjunto de dados favorece a sobreposição exagerada de patches, na extração por SS.

# Conclusão e Trabalhos Futuros

- Extração Randômica Aprimorada demonstrou melhor resultado entre os métodos dinâmicos;
- O método RA com projeção convolucional utilizando o modelo ViT-Base, com mais camadas descongeladas, mostra-se promissor;
- Melhoria do balanceamento através de Redes Generativas (GANs);
- Outros métodos de extração de patches, que seguem a visão humana;
- Treinamento de mais camadas do modelo ViT-Base.

# Referências I



Barcelos, M. R. B., Lima, R. d. C. D., Tomasi, E., Nunes, B. P., Duro, S. M. S., and Facchini, L. A. (2017). Quality of cervical cancer screening in brazil: external assessment of the pmaq. *REV SAUDE PUBL*, 51:67.



Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., and Hounsby, N. (2021). An image is worth 16x16 words: Transformers for image recognition at scale. *ArXiv*, abs/2010.11929.



Kotyan, S. and Vargas, D. V. (2024). Improving robustness for vision transformer with a simple dynamic scanning augmentation. *Neurocomputing*, 565:127000.



N. Diniz, D., T. Rezende, M., G. C. Bianchi, A., M. Carneiro, C., J. S. Luz, E., J. P. Moreira, G., M. Ushizima, D., N. S. de Medeiros, F., and J. F. Souza, M. (2021). A deep learning ensemble method to assist cytopathologists in pap test image classification. *J IMAGING SCI*, 7(7).



Rezende, M. T., Silva, R., Bernardo, F. d. O., Tobias, A. H. G., Oliveira, P. H. C., Machado, T. M., Costa, C. S., Medeiros, F. N. S., Ushizima, D. M., Carneiro, C. M., and Bianchi, A. G. C. (2021). Cric searchable image database as a public platform for conventional pap smear cytology data. *Scientific Data*, 8(1):151.



Wu, B., Xu, C., Dai, X., Wan, A., Zhang, P., Tomizuka, M., Keutzer, K., and Vajda, P. (2021). Visual transformers: Token-based image representation and processing for computer vision. *ArXiv*, abs/2006.03677.



PUC Minas

Obrigado!  
vgiovanini@sga.pucminas.br  
alexeimcmachado@gmail.com

