

Doc Type: Working Group Document
 Title: Feedback on IRGN2578 “Comments on encoding early Chinese organic chemical character in WS2021 and other complex ideographs”¹
 Source: Huáng Jùnliàng (黃俊亮)
 Status: Individual contribution
 Action required: To be considered by the IRG and UTC
 Date: March 12, 2023

In IRGN2578, Eiso Chan listed four base characters of the first type of early organic chemical characters: 𠄎, 𠄏, 𠄐 and 𠄑. Additionally, Chan noted the importance of checking historical usage of numerals greater than 10 when combined with these base characters. To recap, the combined characters are instances of the hydrocarbon nomenclature first proposed by 梁國常 in his work 有機化學命名芻議 [1].

Table 1 contains a list of unencoded hydrocarbon characters based on evidence from Liang’s work.

Table 1: Unencoded hydrocarbon characters in Liang’s work not submitted for WS2021

Character	Page	Evidence
𠄎	73	Saturated hydrocarbons C_nH_{2n+2} 𠄎 (音充)
𠄏	73	Olefines series C_nH_{2n} 𠄏 (音欠)
𠄐	73	Acetylene series C_nH_{2n-2} 𠄐 (音少)
𠄑	73	Benzene ring and Furane ring 𠄑 (音團)
𠄒	73	Propane 𠄒 (讀充三)
𠄓	78	3-chloro-1-propene $CH_2:CH\cdot CH_2Cl$ 𠄓 ¹ 𠄔 ⁸
𠄔	78	allylene C_3H_6 𠄔

¹Sources of this document are available online: <https://github.com/JLHwung/IRGN2578Feedback>.

Character	Page	Evidence
𪛗	73	Butane 𪛗 (讀充四)
𪛘	78	butylene C ₄ H ₈ 𪛘
𪛙	81	pentamethylene diamine NH ₂ (CH ₂) ₅ NH ₂ 𪛙二脰 ^{1:5}
𪛚	73	Furane 𪛚 (讀園五)
		<div style="display: flex; justify-content: space-between;"> <div style="width: 30%;"> <p>furane</p> <p>thiophene</p> <p>pyrrole</p> </div> <div style="width: 40%; text-align: center;"> </div> <div style="width: 20%; text-align: right;"> <p>園 𪛚</p> <p>園 𪛛</p> <p>園 𪛜</p> </div> </div>
	88	
𪛝 ²	81	glyceryl tripalmitate (tripalmitin) C ₃ H ₅ (O·COO ₁₅ H ₃₁) ₃ 三 𪛝酸 𪛞
	81	monopalmitin (HO) ₂ C ₈ H ₅ OCOC ₁₅ H ₃₁ 𪛝酸 𪛞二 𪛟
𪛞	80	monocyclic acid C ₁₈ H ₃₇ COOH 𪛞酸
𪛟	78	heneicosane C ₂₁ H ₄₄ 𪛟
𪛠	78	hexacontane C ₆₀ H ₁₂₂ 𪛠

²To clarify, in the presented evidences, it would be more accurate to use 𪛝 instead of 𪛞, as both tripalmitin and monopalmitin contain hexadecanoate (C₁₆H₃₁O₂⁻) rather than pentadecanoate (C₁₅H₂₉O₂⁻).

If we are to encode 𤝵 (00016), 𤝶 (00017), 𤝷 (01900), 𤝸 (00014) and 𤝹 (00777) one by one, then we should also encode the other 15 characters above in order to digitize Liang's work. In total, we need to encode 20 characters, which is already more than the length of a minimal Unicode block (16 characters), should they be encoded as IDJ sequences proposed in IRGN2578.

Additionally, it's important to note that hydrocarbons can have an arbitrary number of carbon atoms, making such characters an open set. It's possible that Liang's nomenclature was adopted in other historical documents, which may contain additional hydrocarbon characters that we're currently unaware of. In other words, we may have to encode more than 20 hydrocarbon characters in the future.

To summarize, encoding hydrocarbon characters as IDJ sequences can help us save codepoint space and address the open set nature of Liang's nomenclature. As presented in IRGN2578, the IDJ approach potentially solve other encoding issues, too. In light of this I recommend that we carefully consider Chan's approach and inform our encoding strategy.

Acknowledgements

"Hulenkius" helps me polish the font.

References

- [1] 梁國常, "有機化學命名芻議" 北京大學月刊, vol. 7, pp. 71-89, 1920. [Online]. Available: http://read.nlc.cn/allSearch/searchDetail?searchType=all&showType=1&indexName=data_404&fid=01J000340

(End of Document)